

**II Università di Roma, Tor Vergata**  
**Dipartimento d'Ingegneria Civile e Ingegneria Informatica**  
**LM in Ingegneria dell'Informazione e dell'Automazione**  
**Complementi di Probabilità e Statistica - Advanced Statistics**  
**Instructors: Roberto Monte & Massimo Regoli**  
**Solved Problems on Confidence and Prediction Intervals 2021-12-17**

**Problem 1** Assume that the log-returns of a stock in a financial market are Gaussian distributed with unknown mean  $\mu$  and variance  $\sigma^2$ . Let  $X$  be the normal random variable representing the realization of the log-returns and let  $X_1, \dots, X_n$  be a simple random sample of size  $n$  drawn from  $X$ . Assume that  $n = 5$  and the realizations of the sample are

$$x_1 \equiv -1.5, \quad x_2 \equiv -0.5, \quad x_3 \equiv 1.5, \quad x_4 \equiv 2.0, \quad x_5 \equiv 2.5$$

1. Determine a 99% confidence interval for the mean  $\mu$ .
2. Find the confidence for an interval of width 0.1.
3. Determine a 90% confidence interval for the standard deviation  $\sigma$ .

**Solution.**

1. From data we obtain

$$\bar{x}_5 \equiv \frac{1}{5} \sum_{k=1}^5 x_k = 0.8$$

and

$$s_{X,5}^2 = \frac{1}{4} \sum_{k=1}^5 (x_k - \bar{x}_5)^2 = 2.95 \Rightarrow s_{X,5} = 1.72$$

Now, since  $X$  is Gaussian distributed with unknown variance, to determine a  $100(1 - \alpha)\%$  confidence interval for the mean  $\mu$  the statistic to be considered is

$$\frac{\bar{X}_n - \mu}{S_{X,n}/\sqrt{n}} \sim t_{n-1}.$$

The achievement of a  $100(1 - \alpha)\%$  confidence interval requires to use the  $\alpha/2$  upper and lower critical values  $t_{n-1,\alpha/2}^+ \equiv t_{n-1,1-\alpha/2}$  and  $t_{n-1,\alpha/2}^- \equiv t_{n-1,\alpha/2} = -t_{n-1,1-\alpha/2}$  of  $t_{n-1}$  for  $\alpha = 0.01$ , where  $t_{n-1,1-\alpha/2}$  [resp.  $t_{n-1,\alpha/2}$ ] denotes the  $(1 - \alpha/2)$  [resp.  $\alpha/2$ ]-quantile. In fact, we have

$$\begin{aligned}
 -t_{n-1,1-\alpha/2} < \frac{\bar{X}_n - \mu}{S_{X,n}/\sqrt{n}} < t_{n-1,1-\alpha/2} &\Leftrightarrow -\left(\bar{X}_n + t_{n-1,1-\alpha/2} \frac{S_{X,n}}{\sqrt{n}}\right) < -\mu < -\left(\bar{X}_n - t_{n-1,1-\alpha/2} \frac{S_{X,n}}{\sqrt{n}}\right) \\
 &\Leftrightarrow \bar{X}_n - t_{n-1,1-\alpha/2} \frac{S_{X,n}}{\sqrt{n}} < \mu < \bar{X}_n + t_{n-1,1-\alpha/2} \frac{S_{X,n}}{\sqrt{n}},
 \end{aligned}$$

which implies

$$\begin{aligned}
 1 - \alpha &= \mathbf{P} \left( -t_{n-1,1-\alpha/2} < \frac{\bar{X}_n - \mu}{S_{X,n}/\sqrt{n}} < t_{n-1,1-\alpha/2} \right) \\
 &= \mathbf{P} \left( \bar{X}_n - t_{n-1,1-\alpha/2} \frac{S_{X,n}}{\sqrt{n}} < \mu < \bar{X}_n + t_{n-1,1-\alpha/2} \frac{S_{X,n}}{\sqrt{n}} \right).
 \end{aligned}$$

It follows that the desired confidence interval for  $\mu$  is given by the random interval

$$\left( \bar{X}_n - t_{n-1,1-\alpha/2} \frac{S_{X,n}}{\sqrt{n}}, \bar{X}_n + t_{n-1,1-\alpha/2} \frac{S_{X,n}}{\sqrt{n}} \right)$$

A realization of such a confidence interval is then given by

$$\left( \bar{x}_n - t_{n-1, 1-\alpha/2} \frac{s_{X,n}}{\sqrt{n}}, \bar{x}_n + t_{n-1, 1-\alpha/2} \frac{s_{X,n}}{\sqrt{n}} \right).$$

In the case considered, since  $t_{n-1, 1-\alpha/2} \equiv t_{4, 0.995} = 4.60$ ,  $\bar{x}_n \equiv \bar{x}_5 = 0.80$ ,  $s_{X,n} \equiv s_{X,5} = 1.72$ , the realization of the confidence interval becomes

$$(-3.16, 4.76).$$

2. From 1. it is clearly seen the width  $w$  of a  $100(1-\alpha)\%$  confidence interval is given by

$$w = 2t_{n-1, 1-\alpha/2} \frac{S_{X,n}}{\sqrt{n}}.$$

As a consequence, the size  $n$  of the sample which gives a  $100(1-\alpha)\%$  confidence interval of a given width  $w$  is given by the solution of the equation

$$\frac{n}{t_{n-1, 1-\alpha/2}^2} = \left[ 4 \frac{S_{X,n}^2}{w^2} \right] + 1.$$

To determine  $n$  we need to consider as many realizations  $x_1, \dots, x_n$  of the simple random sample  $X_1, \dots, X_n$  such that

$$\frac{n}{t_{n-1, 1-\alpha/2}^2} = \left[ 4 \frac{s_{X,n}^2}{0.01} \right] + 1.$$

3. Again, since  $X$  is Gaussian distributed, to determine a  $100(1-\alpha)\%$  confidence interval for the standard deviation  $\sigma$  the statistic to be considered is

$$\frac{(n-1) S_{X,n}^2}{\sigma^2} \sim \chi_{n-1}^2.$$

Since  $\chi_{n-1}^2$  is not symmetric, the achievement of a  $100(1-\alpha)\%$  confidence interval requires to use the  $\alpha/2$  and the  $1-\alpha/2$  critical value  $\chi_{n-1, \alpha/2}^{2,-} \equiv \chi_{n-1, \alpha/2}^2$  and  $\chi_{n-1, \alpha/2}^{2,+} = \chi_{n-1, 1-\alpha/2}^2$  of  $\chi_{n-1}^2$  for  $\alpha = 0.1$ , where  $\chi_{n-1, \alpha/2}^2$  [resp.  $\chi_{n-1, 1-\alpha/2}^2$ ] is the  $\alpha/2$ -quantile [ $1-\alpha/2$ -quantile] of the  $\chi_{n-1}^2$  distribution. In fact, we have

$$\begin{aligned} \chi_{n-1, \alpha/2}^{2,-} &< \frac{(n-1) S_{X,n}^2}{\sigma^2} < \chi_{n-1, \alpha/2}^{2,+} \Leftrightarrow \frac{1}{\chi_{n-1, \alpha/2}^{2,-}} > \frac{\sigma^2}{(n-1) S_{X,n}^2} > \frac{1}{\chi_{n-1, \alpha/2}^{2,+}} \\ &\Leftrightarrow \frac{(n-1) S_{X,n}^2}{\chi_{n-1, \alpha/2}^{2,+}} < \sigma^2 < \frac{(n-1) S_{X,n}^2}{\chi_{n-1, \alpha/2}^{2,-}}, \end{aligned}$$

which implies

$$1 - \alpha = \mathbf{P} \left( \chi_{n-1, \alpha/2}^{2,-} < \frac{(n-1) S_{X,n}^2}{\sigma^2} < \chi_{n-1, \alpha/2}^{2,+} \right) = \mathbf{P} \left( \frac{(n-1) S_{X,n}^2}{\chi_{n-1, \alpha/2}^{2,+}} < \sigma^2 < \frac{(n-1) S_{X,n}^2}{\chi_{n-1, \alpha/2}^{2,-}} \right).$$

It follows that the desired confidence interval for the variance  $\sigma^2$  is given by

$$\left( \frac{(n-1) S_{X,n}^2}{\chi_{n-1, \alpha/2}^{2,+}} < \sigma^2 < \frac{(n-1) S_{X,n}^2}{\chi_{n-1, \alpha/2}^{2,-}} \right) = \left( \frac{(n-1) S_{X,n}^2}{\chi_{n-1, 1-\alpha/2}^2}, \frac{(n-1) S_{X,n}^2}{\chi_{n-1, \alpha/2}^2} \right).$$

In the case considered, since  $\chi_{n-1,\alpha/2}^2 \equiv \chi_{4,0.5}^2 = 0.71$ ,  $\chi_{n-1,1-\alpha/2}^2 \equiv \chi_{4,0.95}^2 = 9.49$ ,  $\bar{x}_n \equiv \bar{x}_5 = 0.80$ ,  $s_{X,n}^2 \equiv s_{X,5}^2 = 2.95$ , a realization of the confidence interval is given by

$$\left( \frac{4s_{X,n}^2}{\chi_{4,0.95}^2}, \frac{4s_{X,n}^2}{\chi_{4,0.05}^2} \right) = \left( \frac{4 \cdot 2.95}{9.49}, \frac{4 \cdot 2.95}{0.71} \right) = (1.24, 16.62).$$

As a consequence the  $100(1 - 0.1)\%$  confidence interval for the standard deviation  $\sigma$  is

$$(1.11, 4.08).$$

This completes the solution.

**Problem 2** Assume that a library master believes that the mean duration in days of the borrowing period is 20d. However, the library master selects randomly a simple random sample of 100 books in the library and discovers that the sample mean and variance of the borrowing days are 18d and  $8d^2$ , respectively. Determine a 99% confidence interval for the mean duration of the borrowing days to check whether library master's initial guess is correct.

**Solution.** Note that the distribution of the random variable  $X$  representing the duration in days of the borrowing period is unknown. However, are known the sample mean and variance realizations referred to a simple sample of size  $n = 100$ , which may be considered a large sample. In this case the statistic to be considered is given by

$$\frac{\bar{X}_n - \mu}{S_{X,n}/\sqrt{n}},$$

which is approximatively distributed as a standard Gaussian random variable  $Z \sim N(0, 1)$ . The achievement of a  $100(1 - \alpha)\%$  confidence interval requires to use the  $\alpha/2$  upper and lower critical values  $z_{\alpha/2}^+ \equiv z_{1-\alpha/2}$  and  $z_{\alpha/2}^- \equiv z_{\alpha/2} = -z_{1-\alpha/2}$  of  $Z$  for  $\alpha = 0.01$ , where  $z_{1-\alpha/2}$  [resp.  $z_{\alpha/2}$ ] denotes the  $(1 - \alpha/2)$  [resp.  $\alpha/2$ ]-quantile. In fact, we have

$$\begin{aligned} -z_{1-\alpha/2} < \frac{\bar{X}_n - \mu}{S_{X,n}/\sqrt{n}} < z_{1-\alpha/2} &\Leftrightarrow -\left(\bar{X}_n + z_{1-\alpha/2} \frac{S_{X,n}}{\sqrt{n}}\right) < -\mu < -\left(\bar{X}_n - z_{1-\alpha/2} \frac{S_{X,n}}{\sqrt{n}}\right) \\ &\Leftrightarrow \bar{X}_n - z_{1-\alpha/2} \frac{S_{X,n}}{\sqrt{n}} < \mu < \bar{X}_n + z_{1-\alpha/2} \frac{S_{X,n}}{\sqrt{n}}, \end{aligned}$$

whuihc implies

$$\begin{aligned} 1 - \alpha &\approx \mathbf{P}\left(-z_{1-\alpha/2} < \frac{\bar{X}_n - \mu}{S_{X,n}/\sqrt{n}} < z_{1-\alpha/2}\right) \\ &= \mathbf{P}\left(\bar{X}_n - z_{1-\alpha/2} \frac{S_{X,n}}{\sqrt{n}} < \mu < \bar{X}_n + z_{1-\alpha/2} \frac{S_{X,n}}{\sqrt{n}}\right) \end{aligned}$$

It follows that the desired confidence interval for  $\mu$  is given by

$$\left(\bar{X}_n - z_{1-\alpha/2} \frac{S_{X,n}}{\sqrt{n}}, \bar{X}_n + z_{1-\alpha/2} \frac{S_{X,n}}{\sqrt{n}}\right).$$

In the case considered,

$$n = 100, \quad \bar{x}_n \equiv \bar{x}_{100} = 18, \quad s_{X,n} \equiv s_{X,100} = \sqrt{8}, \quad z_{1-\alpha/2} \equiv 2.58.$$

Therefore, a realization of the confidence interval is given by

$$\left( \bar{x}_n - z_{1-\alpha/2} \frac{s_{X,n}}{\sqrt{n}}, \bar{x}_n + z_{1-\alpha/2} \frac{s_{X,n}}{\sqrt{n}} \right) = \left( 18 - 2.58 \cdot \frac{\sqrt{8}}{\sqrt{100}}, 18 + 2.58 \cdot \frac{\sqrt{8}}{\sqrt{100}} \right) = (17.27, 18.73).$$

It follows that library master's initial guess is not supported by data. Note that this problem can be tackled also exploiting the hypothesis test method. In fact, assume as the null hypothesis that library master's assumption is correct, that is  $H_0 : \mu = \mu_0$ , and as the alternative hypothesis that library master's assumption is wrong, that is  $H_0 : \mu \neq \mu_0$ . The same consideration as above on the available information on the random variable  $X$  led to consider the rejection region

$$R = \left\{ \frac{\bar{X}_n - \mu}{s_{X,n}/\sqrt{n}} < z_{\alpha/2} \right\} \cup \left\{ \frac{\bar{X}_n - \mu}{s_{X,n}/\sqrt{n}} > z_{1-\alpha/2} \right\} = \left\{ \frac{\bar{X}_n - \mu}{s_{X,n}/\sqrt{n}} < -2.58 \right\} \cup \left\{ \frac{\bar{X}_n - \mu}{s_{X,n}/\sqrt{n}} > 2.58 \right\},$$

where  $\mu = 20$ . Computing the statistic  $\frac{\bar{X}_n - \mu}{s_{X,n}/\sqrt{n}}$  for the available realization, we obtain

$$\frac{\bar{x}_n - \mu}{s_{X,n}/\sqrt{n}} = \frac{18 - 20}{\sqrt{8}/10} = -7.07 \in R.$$

Hence, the library master's assumption has to be rejected.

**Problem 3** A sample of 60 cars of the same model are tested for gasoline consumption, expressed as litres per 100 kilometers (L/100km). The result of the test yield a mean consumption  $\mu$  of 9.4 L/100km and a standard deviation  $\sigma$  of 1.5 L/100km.

1. Determine a 99% confidence interval for the mean consumption. Is it necessary to make any assumption on the consumption distribution?
2. Assume that the consumption is normally distributed with variance  $\sigma^2 = 2$  L/100km. How large has to be the sample if, with the same confidence, we want the maximum error to be 0.25 L/100km? Apply the method of the confidence interval and the Chebychev inequality.

**Solution.**

**Problem 4** Let  $X$  [resp.  $Y$ ] a Gaussian distributed random variables with (unknown) mean  $\mu_X \in \mathbb{R}$  [resp.  $\mu_Y \in \mathbb{R}$ ] and variance  $\sigma_X^2 > 0$  [resp.  $\sigma_Y^2 > 0$ ]. Assume that  $X$  describes a population before a treatment and  $Y$  describes the same population after a treatment. Let  $X_1, \dots, X_n$  be a simple random sample drawn by  $X$  and let  $Y_1, \dots, Y_n$  be the corresponding sample drawn from  $Y$ . Note that we can still assume that  $Y_1, \dots, Y_n$  is a simple random sample but we cannot assume that the samples  $X_1, \dots, X_n$  and  $Y_1, \dots, Y_n$  are independent. Actually, there is no reason at all to think that the random variables  $X$  and  $Y$  are independent. However, it is still reasonable to assume that the random variable  $D \equiv Y - X$  is Gaussian distributed and that  $D_1 \equiv Y_1 - X_1, \dots, D_n \equiv Y_n - X_n$  is a simple random sample drawn from  $D$ .

1. Given  $\alpha > 0$ , can you build a  $100(1 - \alpha)\%$  confidence interval for the difference  $\mu_Y - \mu_X$ ?
2. Assume to have measured

$$\begin{array}{cccccccc} x_1 = 3.85 & x_2 = 2.82 & x_3 = 3.44 & x_4 = 3.48 & x_5 = 1.92 & x_6 = 4.39 & x_7 = 3.12 \\ y_1 = 5.73 & y_2 = 3.84 & y_3 = 4.78 & y_4 = 4.40 & y_5 = 1.91 & y_6 = 4.98 & y_7 = 4.94 \end{array}.$$

Determine a realization of the 95% confidence interval built above.

**Solution.** We clearly have

$$\mu_D \equiv \mathbf{E}[D] \equiv \mathbf{E}[Y - X] = \mathbf{E}[Y] - \mathbf{E}[X] = \mu_Y - \mu_X.$$

Therefore, introducing the sample mean of size  $n$  drawn from  $D$ , that is

$$\bar{D}_n = \frac{1}{n} \sum_{k=1}^n D_k = \frac{1}{n} \sum_{k=1}^n (Y_k - X_k) = \frac{1}{n} \sum_{k=1}^n Y_k - \frac{1}{n} \sum_{k=1}^n X_k = \bar{Y}_n - \bar{X}_n$$

and the unbiased sample variance of size  $n$  drawn from  $D$ , that is

$$S_n^2(D) = \frac{1}{n-1} \sum_{k=1}^n (D_k - \bar{D}_n)^2,$$

under the assumptions considered, we have that the statistic

$$\frac{\bar{D}_n - \mu_D}{S_n(D) / \sqrt{n}},$$

has the Student distribution with  $n - 1$  degrees of freedom. As a consequence, given  $\alpha > 0$  a  $100(1 - \alpha)\%$  confidence interval for  $\mu_D$  is given by

$$\left( \bar{D}_n - t_{\frac{\alpha}{2}, n-1}^- \frac{S_n(D)}{\sqrt{n}}, \bar{D}_n + t_{\frac{\alpha}{2}, n-1}^+ \frac{S_n(D)}{\sqrt{n}} \right).$$

The realization of such a confidence interval are of the form

$$\left( \bar{d}_n - t_{\frac{\alpha}{2}, n-1} \frac{s_n(D)}{\sqrt{n}}, \bar{d}_n + t_{\frac{\alpha}{2}, n-1} \frac{s_n(D)}{\sqrt{n}} \right),$$

where  $\bar{d}_n$  [resp.  $s_n(D)$ ] is the value taken by the sample mean estimator  $\bar{D}_n$  [resp. unbiased sample standard deviation estimator  $S_n(D)$ ] on the available realizations  $d_1, \dots, d_n$  of the sample  $D_1, \dots, D_n$ . Since in our case  $n = 7$  and  $\alpha \equiv 0.05$ , we have

$$t_{\frac{\alpha}{2}, n-1} \equiv t_{0.025, 6} = 2.45$$

Furthermore,

$$\bar{d}_n \equiv \bar{d}_7 = 1.08 \quad \text{and} \quad s_n(D) \equiv s_7(D) = 0.45$$

Therefore,

$$\bar{d}_n - t_{\frac{\alpha}{2}, n-1} \frac{s_n(D)}{\sqrt{n}} = 1.08 - 2.45 \frac{0.45}{\sqrt{7}} = 0.66$$

and

$$\bar{d}_n + t_{\frac{\alpha}{2}, n-1} \frac{s_n(D)}{\sqrt{n}} = 1.08 + 2.45 \frac{0.45}{\sqrt{7}} = 1.50.$$

Thus the realization of the confidence interval is

$$(0.66, 1.50).$$

**Problem 5** Two independent groups of 50 people, say group A and group B, are affected by a disease. A drug is given to the individuals of group A (treatment group) but not to the individuals of group B (control group). After a week since the administration of the drug a medical check shows that 45 [resp. 30] people of group A [resp. B] got recovered from the disease. Let  $p_A$  [resp.  $p_B$ ] the proportion of individuals of group A [resp. B] who got recovered.

1. Given  $\alpha > 0$ , can you build a  $100(1 - \alpha)\%$  confidence interval for the true value of difference  $p_A - p_B$ ?
2. Determine the realization of the 95% confidence interval built above and use it to comment on the efficacy of the drug.

**Solution.** .

**Problem 6** Let  $X$  be a standard Bernoulli random variable with unknown success parameter  $p$ . Let  $X_1, \dots, X_n$  be a simple random sample of size  $n$  drawn from  $X$  and let  $Z_n \equiv \sum_{k=1}^n X_k$  be the sample sum. It is well known that  $Z_n \sim \text{Bin}(n, p)$ . In addition, when  $n$  is large ( $np \geq 10$  and  $n(1 - p) \geq 10$ ) the sample sum has approximately a normal distribution.

1. Determine a confidence interval for the parameter  $p$  with confidence level approximately  $100(1 - \alpha)\%$ .
2. Determine the size  $n$  of the sample  $X_1, \dots, X_n$  which allows a confidence interval for the parameter  $p$  with confidence level approximately  $100(1 - \alpha)\%$  and width  $w$ , where both  $\alpha$  and  $w$  are given in advance.

**Solution.** .

**Problem 7** Let  $X_1, \dots, X_n, X_{n+1}$  be a simple random sample of size  $n + 1$  drawn from a Gaussian distributed random variable  $X$  with unknown mean  $\mu$  and variance  $\sigma^2$ . Assume that we have observed  $X_1, \dots, X_n$  and we want use the observed values  $x_1, \dots, x_n$  to determine a confidence interval for the prediction of  $X_{n+1}$ . To this goal give detailed answers to the following questions:

1. what is the distribution of the statistic  $\bar{X}_n$ ?
2. what is the distribution of the statistic  $(X_{n+1} - \bar{X}_n) / \sigma \sqrt{1 + 1/n}$ ?
3. what is the distribution of the statistic  $S_n^2 \equiv \frac{1}{n-1} \sum_{k=1}^n (X_k - \bar{X}_n)^2$ ?
4. are the statistics  $X_{n+1} - \bar{X}_n$  and  $S_n^2 \equiv \frac{1}{n-1} \sum_{k=1}^n (X_k - \bar{X}_n)^2$  independent? Why?
5. what is the distribution of the statistic  $(X_{n+1} - \bar{X}_n) / S_n \sqrt{1 + 1/n}$ ?
6. After answering the above questions, build an interval in which the random variable  $X_{n+1}$  takes its values with probability  $\alpha$  and determine the corresponding confidence interval for the prediction of  $X_{n+1}$ . In the end, assume that  $n = 7$  and we have

$$x_1 = 7005, \quad x_2 = 7432, \quad x_3 = 7420, \quad x_4 = 6822, \quad x_5 = 6752, \quad x_6 = 5333, \quad x_7 = 6552.$$

compute the 95% confidence interval for the prediction of  $X_8$ .

**Solution.** .

**Problem 8** A test on the reaction time measured in seconds to a sudden emergency has lead to the following results in 10 people:

$$\begin{aligned} t_1 = 0.77, \quad t_2 = 0.75, \quad t_3 = 0.70, \quad t_4 = 0.72, \quad t_5 = 0.70, \\ t_6 = 0.69, \quad t_7 = 0.67, \quad t_8 = 0.79, \quad t_9 = 0.64, \quad t_{10} = 0.72. \end{aligned}$$

Assume that the reaction time can be modelled by a normal random variable  $T \sim N(\mu, \sigma^2)$  with unknown  $\mu$  and  $\sigma^2$ .

1. Compute the sample mean and variance referred to the above sample.
2. Find the confidence interval for  $\mu$  [resp.  $\sigma^2$ ] at the confidence level 95%.
3. What would the confidence interval for  $\mu$  be if the variance were known and we had  $\sigma^2 = 0.0025$ ?
4. What should the size of the sample be to achieve a precision of 10?

**Solution.** .

**Problem 9** Assume we need to measure the same trait  $X$  in two different population and the results of our measurement for a sample of size  $n_1 = 15$  [resp.  $n_2 = 20$ ] of the first [resp. second] population gives a value  $\bar{x}_{n_1} = 24.0$  [resp.  $\bar{x}_{n_2} = 26.0$ ] of the sample mean  $\mu_X$  [resp.  $\mu_Y$ ] of the characteristic under investigation, with a sample variance  $s_{n_1}^2 = 4.5$  [resp.  $s_{n_2}^2 = 5.0$ ]. Assume that the trait is normally distributed with the same unknown variance  $\sigma_X^2$ . Compute a confidence interval for the value of the difference  $\mu_X - \mu_Y$  with a 95% confidence level.

**Solution.** .