

CSE Exercises - Week 6

- ① Let X be a random variable such that $E(X)$ and $E(X^2)$ exist and are finite. The WLLN (Proposition 50 on page 151) justifies the use of the sample mean, \bar{X}_n , as an estimator for $E(X)$. In this exercise, we investigate how well \bar{X}_n estimates $E(X)$ by constructing "confidence bounds" for the estimation error.

The error in using \bar{X}_n to estimate $E(X)$ is given by $|\bar{X}_n - E(X)|$, where the absolute value is taken because we are interested in the magnitude of the error and disregarding the sign. For $0 < p < 1$, we say that B_p is a $100p\%$ confidence bound when

$$P(|\bar{X}_n - E(X)| \leq B_p) = p,$$

i.e. the estimation error is bounded by B_p with probability p or, equivalently, we can be $100p\%$ confident that the estimation error is at most B_p . For example, when $p = 0.95$, $B_{0.95}$ is a 95% confidence bound.

Download the file, `glass.mat`, to your Matlab working directory. It contains a vector, X , containing the measured refractive

indices of 214 glass specimens. You can load this vector into your Matlab workspace by executing the following command in the Matlab command window:

```
>> load glass
```

Assume that the measured refractive indices are continuous and independent and come from a common distribution with $E(X) < \infty$ and $E(X^2) < \infty$. Our goal is to estimate $E(X)$ using \bar{X}_n and obtain 95% confidence bounds for the estimation error.

(a) Find the sample mean \bar{X}_n .

(b) Use Chebychev's inequality (equation (8.5) on page 150) to find an approximate bound for the estimation error, with at least 95% confidence. The bound is approximate because you will have to estimate $V(X)$ using the sample variance S_n^2 .

(c) Use the form of the CLT given in Proposition 52 (page 153) to find an approximate 95% confidence bound for the estimation error. Explain why the bound obtained here is approximate.

(d) Compare the confidence bounds obtained in parts (b) and (c). What can you conclude?

- ② Let X be a random variable with $E(X) < \infty$ and $E(X^2) < \infty$. In this exercise, we will discover, through computational experiments, that the WLLN and CLT continue to hold for any "measurable" function g of X such that $E[g(X)] < \infty$ and $E[g(X)^2] < \infty$. A function g is "measurable" when $g(X)$ remains a random variable.

Now let X_1, \dots, X_n be independent and identically distributed $\text{Normal}(0,1)$ random variables.

(a) Consider $Y_i = g(X_i) = X_i^2$ for $i=1, \dots, n$. It can be shown that Y_1, \dots, Y_n are independent and identically distributed with a chi-squared distribution with parameter 1. Moreover, $E(Y_i) = 1$ and $V(Y_i) = 2$.

(i) For $n=100$, generate $X_1, \dots, X_n \sim \text{Normal}(0,1)$.

(ii) Compute $Y_i = X_i^2$, $i=1, \dots, n$.

(iii) Compute \bar{Y}_n as an estimate of $E(Y_i)$.

(iv) Use Proposition 52 (page 153) to find an approximate 95% confidence bound for the estimation error.

(v) What can you conclude from the results in (iii) and (iv)?

(b) Consider $Y_i = g(X_i) = e^{X_i}$, $i=1, \dots, n$.
 It can be shown that Y_1, \dots, Y_n are independent and identically distributed with a Lognormal(0,1) distribution (recall Model 19 on page 141). Thus, $E(Y_i) = \sqrt{e}$ and $V(Y_i) = e(e-1)$.

Repeat (i) - (v) in part (a) but with $Y_i = e^{X_i}$ in (ii).

Solutions

① (a) $\bar{X}_n = 0.3654$.

(b) Recall that

$$E(\bar{X}_n) = E(X) \quad \text{and} \quad V(\bar{X}_n) = V(X)/n .$$

Replace X by \bar{X}_n in equation (8.5) :

$$P(|\bar{X}_n - E(\bar{X}_n)| \geq \varepsilon) \leq \frac{V(\bar{X}_n)}{\varepsilon^2}$$

$$\Leftrightarrow P(|\bar{X}_n - E(X)| \geq \varepsilon) \leq \frac{V(X)}{n\varepsilon^2}$$

$$\Leftrightarrow P(|\bar{X}_n - E(X)| \leq \varepsilon)$$

$$= 1 - P(|\bar{X}_n - E(X)| \geq \varepsilon)$$

$$\geq 1 - \frac{V(X)}{n\varepsilon^2} .$$

We want the bound ε to hold with at least 95% confidence, i.e. need

$$1 - \frac{V(X)}{n\varepsilon^2} = 0.95$$

$$\Leftrightarrow V(X)/n\varepsilon^2 = 0.05$$

$$\Leftrightarrow \varepsilon = \sqrt{\frac{V(X)}{0.05n}} \approx \sqrt{\frac{S_n^2}{0.05n}}$$

For the data, $S_n^2 = 9.2225$ and $n = 214$ and so

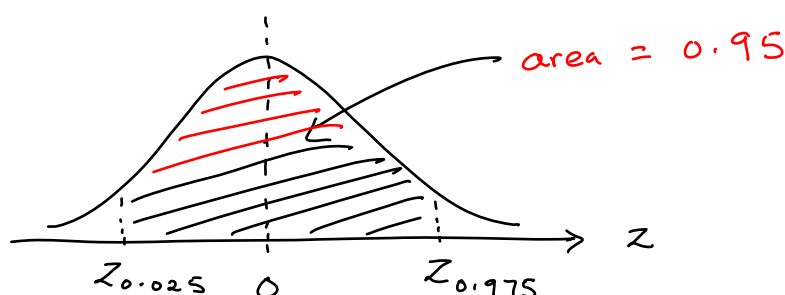
$$\varepsilon \approx 0.9284 .$$

Hence, using Chebychev's inequality, we can be at least 95% confident that the estimation error is at most 0.9284 (approximately).

(c) Since the distribution of $\sqrt{n}(\bar{X}_n - E(X))/S_n$ is approximately Normal(0,1), the following holds:

$$P(Z_{0.025} \leq \frac{\sqrt{n}(\bar{X}_n - E(X))}{S_n} \leq Z_{0.975}) \approx 0.95$$

where for $0 < p < 1$, Z_p denotes the p -quantile of the Normal(0,1) distribution.



Note that by the symmetry of the Normal(0,1) PDF about 0, $Z_{0.025} = -Z_{0.975}$. Thus, we can write

$$P(-Z_{0.975} \leq \frac{\sqrt{n}(\bar{X}_n - E(X))}{S_n} \leq Z_{0.975}) \approx 0.95$$

$$\Leftrightarrow P\left(\frac{\sqrt{n}|\bar{X}_n - E(X)|}{S_n} \leq Z_{0.975}\right) \approx 0.95$$

$$\Leftrightarrow P\left(|\bar{X}_n - E(X)| \leq \frac{Z_{0.975} S_n}{\sqrt{n}}\right) \approx 0.95.$$

Note that the probability is approximately 0.95 and not equal to 0.95 because, for finite n ,

the distribution of $\frac{\sqrt{n}(\bar{X}_n - E(X))}{S_n}$ is approximately Normal $(0,1)$ and not exactly Normal $(0,1)$ to begin with.

Therefore, an approximate 95% confidence bound for the estimation error is given by $\frac{Z_{0.975} S_n}{\sqrt{n}}$

Putting $Z_{0.975} = 1.96$, $S_n^2 = 9.2225$ and $n = 214$, the resulting approximate 95% confidence bound is 0.4069.

Thus, we can be approximately 95% confident that the estimation error is at most 0.4069.

(d) The bound provided by Chebychev's inequality is bigger than the bound provided by the CLT. However, note that whilst the CLT bound gives a confidence of 95% approximately, the Chebychev bound gives a confidence of at least 95% approximately.

In practice, the CLT bound is preferred because it is tighter (smaller) than the Chebychev bound.

(2) (a) (iii) $\bar{Y}_n = 1.1928$.

(iv) Approximate 95% CLT confidence bound is

$$\frac{Z_{0.975} S_n}{\sqrt{n}} = \frac{1.96 \times 2.0926}{10} = 0.4023.$$

(v) Since we know $E(Y_i) = 1$, the estimation error in this case is

$$|\bar{Y}_n - E(Y_i)| = |1.1928 - 1| = 0.1928$$

which is within the 95% confidence bound obtained in (iv).

(b) (iii) $\bar{Y}_n = 1.9131$.

(iv) Approximate 95% CLT confidence bound is

$$\frac{Z_{0.975} S_n}{\sqrt{n}} = \frac{1.96 \times 2.4327}{10} = 0.4768.$$

(v) Here, $E(Y_i) = \sqrt{e} \approx 1.6487$, and so the estimation error is

$$|1.9131 - 1.6487| = 0.2644,$$

which is, once again, within the bound obtained in (iv).