

CS-6015  
Linear Algebra and Random Processes  
Programming Assignment 2

Group No : 12  
Pritha Ganguly (CS17S013)  
Sweta Kumari (BT17D019)  
Indian Institute of Technology Madras

November 6, 2017

Course Instructor : Dr. Prashanth L.A.

## Problem-1

Let  $X_1, \dots, X_N$  denote the sequence of i.i.d random variables, each with mean  $\mu$  and variance  $\sigma^2$ . Let  $\bar{X}_N$  denote the sample mean, i.e

$$\bar{X}_N = \frac{1}{N} \sum_{i=1}^N X_i$$

Therefore,

$$\begin{aligned} E(\bar{X}_N) &= E\left(\frac{1}{N} \sum_{i=1}^N X_i\right) \\ E(\bar{X}_N) &= \frac{1}{N} E(X_1 + X_2 \dots + X_N) \end{aligned}$$

Using linearity of expectation,

$$\begin{aligned} E(\bar{X}_N) &= \frac{1}{N} (E(X_1) + E(X_2) \dots + E(X_N)) \\ &= \frac{1}{N} (N\mu) \\ &= \mu \end{aligned}$$

Similarly,

$$\begin{aligned} Var(\bar{X}_N) &= Var\left(\frac{1}{N} \sum_{i=1}^N X_i\right) \\ Var(\bar{X}_N) &= \left(\frac{1}{N}\right)^2 Var(X_1 + X_2 \dots + X_N) \end{aligned}$$

As all  $X_i$ 's are i.i.d, variance of their sum is the sum of their variances,

$$\begin{aligned} Var(\bar{X}_N) &= \left(\frac{1}{N}\right)^2 (Var(X_1) + Var(X_2) \dots + Var(X_N)) \\ Var(\bar{X}_N) &= \left(\frac{1}{N}\right)^2 (N\sigma^2) \\ &= \frac{\sigma^2}{N} \end{aligned}$$

## Problem-2

Using the theorem of Hoeffding inequality for bounded random variables, we have the following,

$$P(\bar{X}_N - \mu \geq \epsilon) \leq \exp\left(-\frac{2N\epsilon^2}{(b-a)^2}\right) \quad (1)$$

$$P(\bar{X}_N - \mu \leq -\epsilon) \leq \exp\left(-\frac{2N\epsilon^2}{(b-a)^2}\right) \quad (2)$$

We have to show that, for  $\lambda \in (0, 1)$ ,  $\epsilon' > 0$

$$P(\mu \in [\bar{X}_N - \epsilon', \bar{X}_N + \epsilon']) \geq 1 - \lambda$$

Equivalently, the complementary event will be:

$$P(\mu \notin [\bar{X}_N - \epsilon', \bar{X}_N + \epsilon']) \leq \lambda$$

which means,

$$\begin{aligned} P(\mu \leq \bar{X}_N - \epsilon'), P(\mu \geq \bar{X}_N + \epsilon') &\leq \lambda \\ P(\epsilon' \leq \bar{X}_N - \mu), P(-\epsilon' \geq \bar{X}_N - \mu) &\leq \lambda \end{aligned}$$

From equation (1) and (2), we can write,

$$2\exp\left(-\frac{2N\epsilon^2}{(b-a)^2}\right) \leq \lambda$$

Taking log on both sides,

$$\begin{aligned} \log(2) - \frac{2N\epsilon^2}{(b-a)^2} &\leq \log(\lambda) \\ (b-a)^2 \log\left(\frac{2}{\lambda}\right) &\leq 2N\epsilon^2 \\ \sqrt{\frac{(b-a)^2 \log\left(\frac{2}{\lambda}\right)}{2N}} &\leq \epsilon \end{aligned}$$

Therefore,

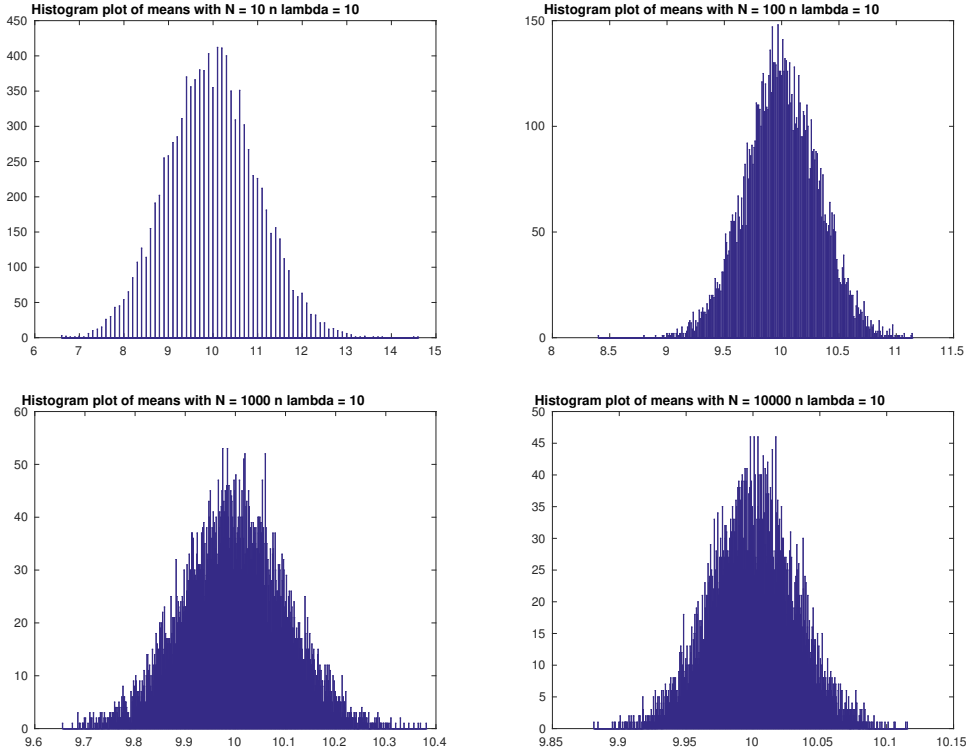
$$\epsilon \geq \sqrt{\frac{(b-a)^2 \log\left(\frac{2}{\lambda}\right)}{2N}}$$

### Problem-3

a) Yes, the sample mean is close to the true mean as we proved in problem-1. Let  $X_1, \dots, X_N$  denote the sequence of i.i.d random variables, each with mean  $\mu$  and variance  $\sigma^2$ . Let  $\bar{X}_N$  denote the sample mean, i.e

$$\bar{X}_N = \frac{1}{N} \sum_{i=1}^N X_i$$

$$\begin{aligned} E(\bar{X}_N) &= \frac{1}{N} (E(X_1) + E(X_2) + \dots + E(X_N)) \\ &= \frac{1}{N} (N\mu) \\ &= \mu \end{aligned}$$



**Figure 1:** Sample mean histograms for  $N = 10, 100, 1000, 10000$ .

The histogram of sample means for different values on  $N$  is given in Figure 1.

**b)** The following table is showing the times was the sample mean in the interval  $[9.99; 10.01]$  and  $[9.9; 10.1]$  for  $N=10, 100, 1000, 10000$ .

N	10	100	1000	10000
interval of 9.99 and 10.01	355	391	822	2456
interval of 9.9 and 10.1	1170	2655	6818	9989

**c)** The following table is showing the times did the true mean fall outside the confidence interval for  $N=10, 100, 1000, 10000$ .

N	10	100	1000	10000
Falls outside the interval	513	508	529	457

**d)** Theorem 1 is only applicable for bounded random variables but Poisson distribution is not bounded hence we cannot apply the Hoeffding bound on it. We can apply the Hoeffding bound by approximating Poisson distribution with parameter  $\lambda$  by a Binomial distribution with parameters  $n$  and  $\frac{\lambda}{n}$ . A binomial distribution is the sum of independent and identically distributed (i.i.d) Bernoulli random variables. Let  $X_1, \dots, X_n$  denote the sequence of i.i.d Bernoulli( $\frac{\lambda}{n}$ ) random variables with each  $\bar{X}_i \in [0, 1]$  for all  $i = 1, \dots, n$ . Therefore, Hoeffding inequality says that,

$$P(|\bar{X}_N - \lambda| \geq \epsilon) \leq 2\exp(-2n\epsilon^2)$$

or,

$$P(|\bar{X}_N - \lambda| < \epsilon) \geq 1 - 2\exp(-2n\epsilon^2)$$

For 95% confidence interval, we must have,

$$1 - 2\exp(-2n\epsilon^2) = 1 - 0.05 = 0.95$$

$$\epsilon = \sqrt{\frac{\log(40)}{2n}}$$

Therefore, 95% confidence interval for  $p = \lambda$  can be written as,

$$P(\bar{X}_N - \sqrt{\frac{\log(40)}{2n}} < \lambda < \bar{X}_N + \sqrt{\frac{\log(40)}{2n}}) \geq 0.95$$

Comparing the above with experimental confidence interval, we get the following table :

N	10	100	1000	10000
Sample mean	9.8	10.36	9.8899	10.03
CI(expt)	[7.40,12.19]	[8.02,12.51]	[8.03,11.95]	[8.08,11.55]
CI(Hoeffding)	[9.37,10.23]	[10.22,10.49]	[9.83,9.92]	[10.017,10.043]

e) The following table is showing the value of N for accuracy 0.1 and 0.01.

Accuracy	0.1	0.01
N	62	611

When the accuracy is 0.001, the sample size is approximately 6100. So, we can generalize and say that on decreasing the accuracy value by one decimal place, the sample size increases 10 times.

## Problem-4

a) Given a random variable  $X$  that takes values  $\pm 1, \pm 2, \dots$  with p.m.f  $f$  defined as:

$$\bar{f}_k = \frac{A}{k^2}$$

for  $k = \pm 1, \pm 2, \dots$  As we know that sum of the total probability is

$$\sum_k f_k = 1$$

Therefore,

$$\sum_k \frac{A}{k^2} = 1$$

Now we can write it as:

$$A\left[\frac{1}{(-1)^2} + \frac{1}{(+1)^2} + \frac{1}{(-2)^2} + \frac{1}{(+2)^2} + \dots\right] = 1$$

$$A\left[\frac{2}{(1)^2} + \frac{2}{(2)^2} + \dots\right] = 1$$

$$2A\left[\frac{1}{(1)^2} + \frac{1}{(2)^2} + \dots\right] = 1$$

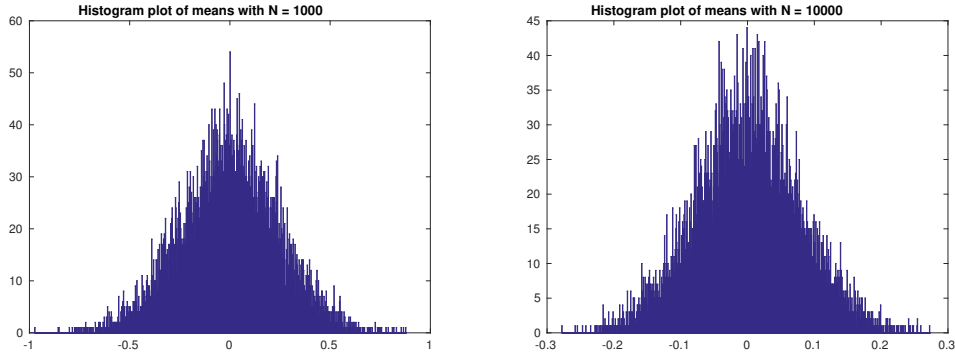
$$2A\frac{\pi^2}{6} = 1$$

$$A = \frac{6}{2\pi^2}$$

$$A = \frac{3}{\pi^2}$$

Therefore for  $A = \frac{3}{\pi^2}$ ,  $f$  would be a valid p.m.f.

**b)** The histogram of sample means for different values on N is given in Figure 2. The sample mean does not stay close to a particular value, it varies. Given



**Figure 2:** Sample mean histograms for  $N = 1000, 10000$ .

the PMF, the expectation of  $X$ ,

$$\sum k \cdot f_k = A \sum_{k \neq 0} \frac{1}{k}$$

does not converge absolutely. Therefore, the sample mean does not converge to a particular value hence confidence intervals cannot be defined.