

# Statistical Inference: Assignment #2

Reyhaneh Ahani  
reyhaneh.ahani@gmail.com

University of Tehran — December 21, 2023

## Introduction

Initially, all the theoretical questions are addressed, followed by the programming section that includes output figures and interpretations of the results.

## 1 Theoretical Section

Handwritten mathematical derivations for a binomial distribution problem. The work is divided into two parts, (a) and (b).

Part (a):

- Given:  $p(X=4) = 1/4$ ,  $p(X \neq 4) = 3/4$ ,  $n=60$ .
- Expected value:  $E(p(X=4)) = \frac{60}{4}$ .
- Variance:  $Var = 60 \times \frac{1}{4} \times \frac{3}{4} = \frac{45}{2}$ .
- Standard deviation:  $\sigma = \sqrt{\frac{45}{2}}$ .
- Probability calculation for  $X=10$ :  
$$P(X=10) = P\left(\frac{X - \frac{60}{4}}{\sqrt{\frac{45}{2}}} = z = \frac{10 - \frac{60}{4}}{\sqrt{\frac{45}{2}}}\right) \rightarrow \Phi(-1.82) = 1 - \Phi(1.82) = 0.0344$$
- Probability calculation for  $X=50$ :  
$$P(X=50) = P\left(\frac{X - \frac{60}{4}}{\sqrt{\frac{45}{2}}} = z = \frac{50 - \frac{60}{4}}{\sqrt{\frac{45}{2}}}\right) \rightarrow \Phi(1.82) = 0.9656$$
- Distance between  $(5, 10) \Rightarrow 0.9656 - 0.0344 = 0.9312$ .

Part (b):

- Expected value of sum of dice in 100 trials:  $E(\text{Sum of dice in 100}) = \frac{100 \times 3.5}{1} = 350$ .
- Variance:  $Var = \frac{100 \times 3.5}{1} = 350$ .
- Standard deviation:  $\sigma = \sqrt{350}$ .
- Probability calculation for  $X=700$ :  
$$P(X=700) = P\left(z = \frac{700 - 350}{\sqrt{350}}\right) = 0.0011$$

Figure 1: Problem 0

1.  $R = \text{proportion of yes}$

a)  $r = (rp-1)q + (1-p)$

we know that  $\rightarrow P(\text{yes}) = P(\text{yes}|q=1) \times P(q=1) + P(\text{yes}|q=r) \times P(q=r)$

$\rightarrow P(\text{yes}) = \underbrace{P(\text{yes}|\text{have A}) \times p}_{\text{as the question said } q} + \underbrace{P(\text{yes}|\text{not have A}) \times (1-p)}_{\rightarrow 1-q}$

$\rightarrow \underline{P(\text{yes})} = qp + (1-q)(1-p) = qp + 1 - p - q + qp = \underline{2qp - q - p + 1} \quad (*)$

b)  $r$  were know, determine  $q$ ?

$\underline{r} = P(\text{yes}) = 2qp - q - p + 1 \rightarrow r = q(rp-1) - p + 1 \rightarrow \underline{\frac{r-1+p}{rp-1} = q}$

c)  $E(R) = r \rightarrow \text{estimation of } q$ ?

$E(R) = E(2qp - q - p + 1) = 2pE(q) - p + 1 - E(q) \Rightarrow rpq - q - p + 1$

From  $(*) \leftarrow q(rp-1) - p + 1 = r$

Figure 2: Problem 1) a

estimation of  $q \rightarrow \hat{q} = \frac{r-1+p}{rp-1} \Rightarrow E(\hat{q}) = E\left(\frac{r-1+p}{rp-1}\right) = \frac{r-1+p}{rp-1}$  unbiased checking

we know that  $\Rightarrow r = 2qp - p + 1 - q \quad (*)$

$E(\hat{q}) = \frac{2qp - p + 1 - q - 1 + p}{rp-1} = \frac{q(rp-1)}{(rp-1)} \Rightarrow E(\hat{q}) = q = E(q) \checkmark \text{ unbiased}$

Figure 3: Problem 1) b

$$e) \text{Var}(q) = \left( \frac{\partial q}{\partial R} \right)^T \text{Var}(R) ; \quad \frac{\partial q}{\partial R} = \frac{1}{r^p - 1} \Rightarrow \left( \frac{1}{r^p - 1} \right)^T \times \frac{r(1-r)}{n}$$

$$d) \text{Var}(R) = \frac{E(R^2) - E(R)^2}{n} = \frac{E(R^2) - r^2}{n} = \frac{E(RU - R) - r^2}{n} = \frac{E(R \cdot R^2) - r^2}{n}$$

$$\frac{E(R) - E(R^2) - r^2}{n} = \frac{r - E(R^2) - r^2}{n} = \frac{r(1-r)}{n}$$

Figure 4: Problem 1) c

First  $\rightarrow X^T X^{-1} r$   
Second  $\rightarrow X^T X^{-1}$

$90\% \text{ CI} \Leftrightarrow z = 1.94$

SE checking  $\left\{ \begin{array}{l} SE(y, r) = \sqrt{\frac{1.94^2 \times 94}{9450}} = 0.14\% \\ SE(\sigma^2) = \sqrt{\frac{12 \times 9}{9450}} = 10\% \end{array} \right.$

SE  $\left\{ \begin{array}{l} SE_1 = 1.94 \times \frac{1.94 \times 1.94 \times 1.94 \times 1.94}{(1.94 \times 1.94)^2} = 17.21 \text{ NVF} \approx 17.21 \\ SE_2 = \frac{1.94 \times 1.94 \times 1.94 \times 1.94}{(1.94 \times 1.94)^2} = 0.004125 \approx 0.004 \end{array} \right.$

SE checking  $\left\{ \begin{array}{l} SE(y, r) = \sqrt{\frac{1.94^2 \times 94}{12251}} = 0.001 = 0.1\% < 1\% \\ SE(\sigma^2) = \sqrt{\frac{12 \times 9}{12251}} = 0.002 = 0.2\% < 1\% \end{array} \right.$

Figure 5: Problem 2

v) a  $\rightarrow$  wrong  $\Rightarrow$  The CLT describes the behavior of the distribution of sample means as the sample size increases.

b  $\rightarrow$  correct  $\checkmark$

c  $\rightarrow$  wrong  $\Rightarrow$  the CLT states that, as the sample size increases, the sampling distribution of the sample means will closely approximate the normal distribution.

d  $\rightarrow$  correct  $\checkmark$

e  $\rightarrow$  wrong  $\Rightarrow$  for a given standard error, higher confidence levels produce wider confidence intervals.

f  $\rightarrow$  wrong  $\Rightarrow$  we are 90% confident that the true population mean falls between  $\bar{x} - 1.645$  and  $\bar{x} + 1.645$ .

g  $\rightarrow$  correct

Figure 6: Problem 3) a

h  $\rightarrow$  wrong  $\Rightarrow$  we are 90% confident that the true population mean age of freshman college students in the state falls within the range of 18.2 and 19.0 years.

i  $\Rightarrow$  correct

Figure 7: Problem 3) b

$$F) X - Y = W$$

$$W \sim (\mu_x - \mu_y, \sigma_x^2 + \sigma_y^2) \rightarrow Z_{\text{score}}(W) = \frac{(\bar{X} - \bar{Y}) - (\mu_x - \mu_y)}{\sqrt{\frac{\sigma_x^2}{n_x} + \frac{\sigma_y^2}{n_y}}} \sim N(0,1)$$

$$\text{marginal error} = z_{\frac{\alpha}{2}} \times \sqrt{\frac{\sigma_x^2}{n_x} + \frac{\sigma_y^2}{n_y}}, \text{ confidence interval is } 90\%, \boxed{z = 1.64}$$

$$r = 1.64 \times \frac{12.12}{\sqrt{n}} \rightarrow \boxed{n \geq 195}$$

Figure 8: Problem 4

a)  $H_0: \mu_1 = \mu_2$  prove:  $H_0$  is rejected only if the confidence interval for the difference of the means doesn't include zero.

$H_0: \mu_1 = \mu_2 \rightarrow$  it has a meaning that there is no difference between  $\mu_1, \mu_2$

$H_1: \mu_1 \neq \mu_2 \rightarrow$  it has a meaning that there is a significant difference between  $\mu_1, \mu_2$

\* since confidence interval provides a range of values within the true difference of the population means is likely to fall, so, if confidence interval doesn't contain zero init, it means that  $\mu_1 \neq \mu_2$  or

$\mu_1 \neq \mu_2 \rightarrow$  which means that there is a significant difference between  $\mu_1, \mu_2$  and it results that  $H_0$  has to be rejected.

Figure 9: Problem 5

4. under what condition for  $c_i$ ,  $\bar{X}_c$  is an unbiased estimator of population mean?

$$\text{unbiased checking} \rightarrow E(\bar{X}_c) = \mu, \quad \bar{X}_c = \sum_{i=1}^n c_i X_i$$

$$E\left(\sum_{i=1}^n c_i X_i\right) = \left(\sum_{i=1}^n c_i E(X_i)\right) = \mu \rightarrow \mu \sum_{i=1}^n c_i = \mu \rightarrow \boxed{\sum_{i=1}^n c_i = 1}$$

condition  $\Rightarrow$  if the summation of  $c_i$ 's is equal to 1, so the  $\bar{X}_c$  will be an unbiased estimation of population mean.

Figure 10: Problem 6

## 2 Programming Section

### 2.1 Problem 7

#### 2.1.1 Histogram of the population values for cancer mortality

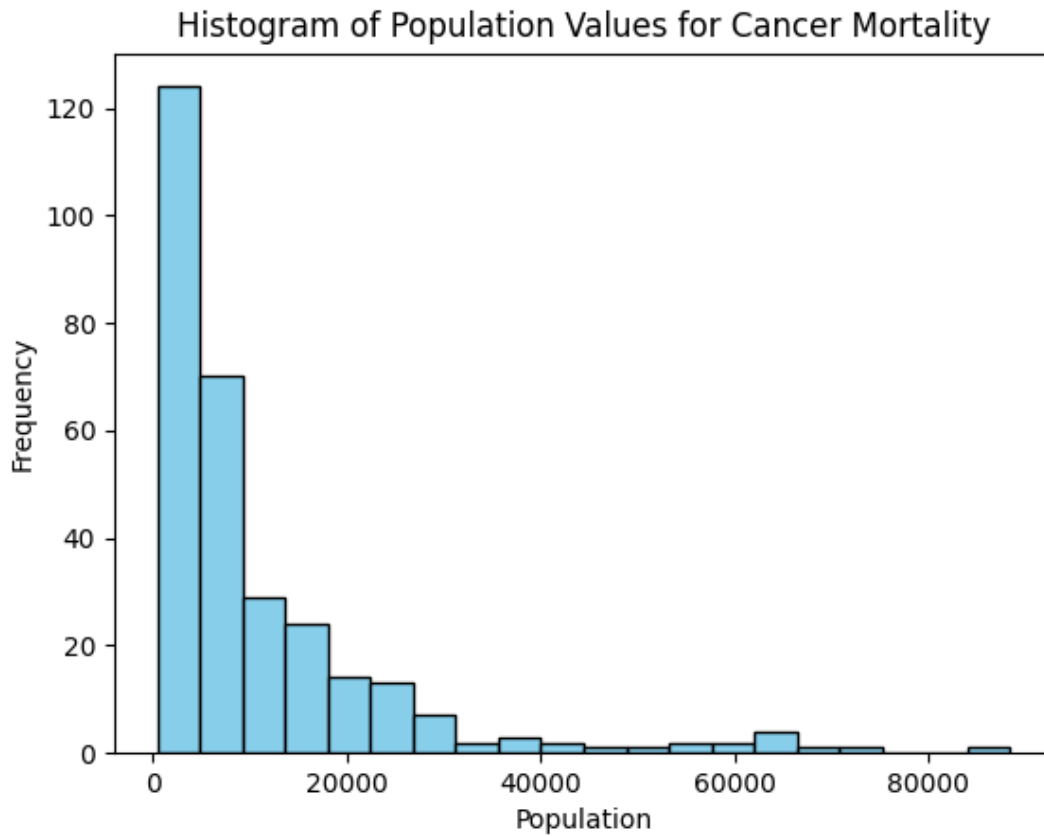


Figure 11: Caption

#### 2.1.2 Population metrics and total cancer mortality

Parameter	Value
Total Cancer Mortality	11997
Population Mean	11288.056478405315
Population Variance	189888678.03346625
Population Standard Deviation	13780.01008829334

Table 1: Population statistics and total cancer mortality

### 2.1.3 Sampling mean distribution from a set of 25 observations

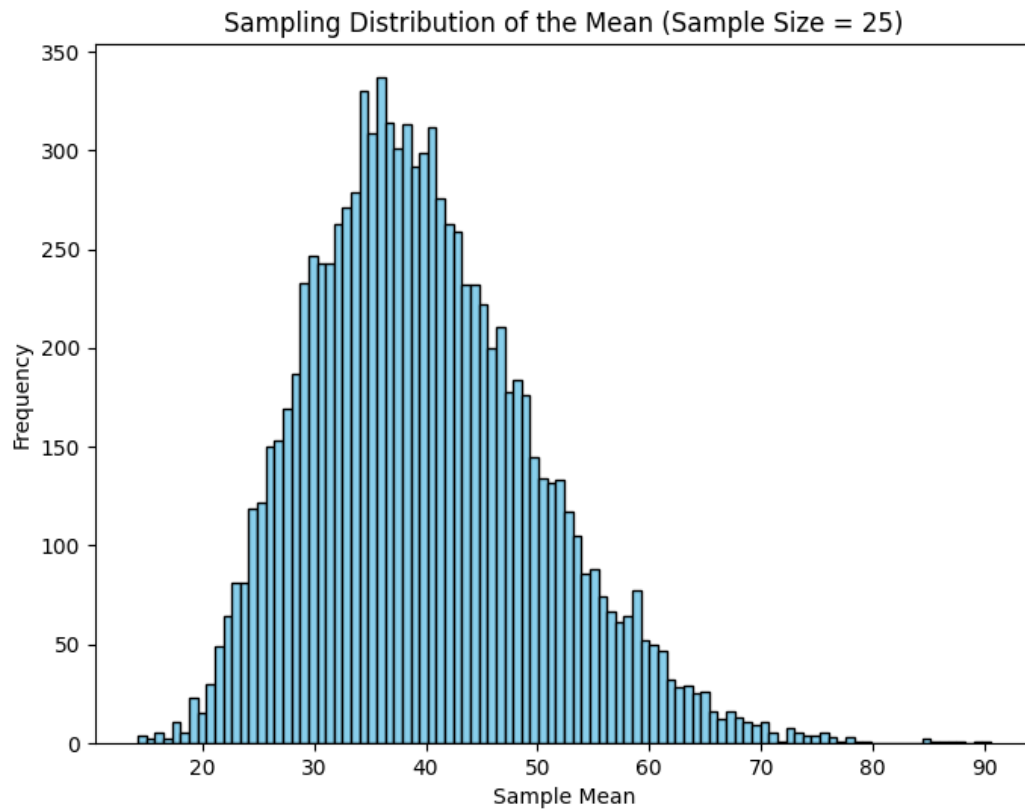


Figure 12: Histogram of sampling mean distribution from a set of 25 observations for  $N = 10000$

### 2.1.4 Estimate the mean and total cancer mortality from a set of 25 observations

Statistic	Value
Sample Mean	30.92
Sample Total Mortality	773
Actual Mortalities Mean	39.86
Actual Total Mortality	11,997
Estimated Mortalities Mean	30.92
Estimated Total Mortality	9,306.92

Table 2: Estimate the mean and total cancer mortality for  $N = 25$

### 2.1.5 Estimate the population variance and standard deviation from a set of 25 observations

Statistic	Value
Estimated Population Var	108141607.35
Estimated Population Var	10399.115

Table 3: Estimate the population variance and standard deviation for  $N = 25$

### 2.1.6 95% confidence intervals

Parameter	Value
Mean Sample Population	8577.16
Mean Population	11288.056
Standard Error	794.266
Margin of Error	1563.039
95% Confidence Intervals	(9725.016, 12851.096)
Intervals Cover Population Values	False

Table 4: 95% confidence intervals statistic Summary for N = 25

### 2.1.7 Repeat parts (d) through (f) for a sample of size 100

1. part d:

Statistic	Value
Sample Mean	37.35
Sample Total Mortality	3735
Actual Mortalities Mean	39.86
Actual Total Mortality	11,997
Estimated Mortalities Mean	37.35
Estimated Total Mortality	11242.349999999999

Table 5: Estimate the mean and total cancer mortality for N = 100

2. part e:

Statistic	Value
Estimated Population Var	179159221.283
Estimated Population Var	13385.037

Table 6: Estimate the population variance and standard deviation for N = 100

3. part f:

Parameter	Value
Mean Sample Population	10849.02
Mean Population	11288.056
Standard Error	794.266
Margin of Error	1563.039
95% Confidence Intervals	(9725.016, 12851.096)
Intervals Cover Population Values	True

Table 7: 95% confidence intervals statistic Summary for N = 25

### 2.1.8 Assessing the Effectiveness of Ratio Estimators in Enhancing Cancer Mortality Estimates Using Total County Population Size

Yes, using a ratio estimator based on the total population of each county could be effective in improving cancer mortality estimates. By incorporating the total population information, the ratio estimator helps to account for variations in population size across different counties. This is important because areas with larger populations may naturally have more cancer cases and mortalities, even if the mortality rate per capita is relatively low. The ratio estimator helps normalize the mortality data by considering the population size, providing more accurate and comparable estimates of cancer mortality rates across counties.



### 2.1.9 Sampling mean distribution from a set of 25 observations for ratio estimators

Part (c) focuses on the mean of the mortality values themselves, while part (i) centers around the ratio of the mean mortality to the mean population size. Information about the population size is incorporated by the ratio estimator, potentially resulting in distinct insights or more efficient estimators, depending on the relationship between mortality and population size. The histograms produced in (c) and (i) mirror these distinct perspectives.

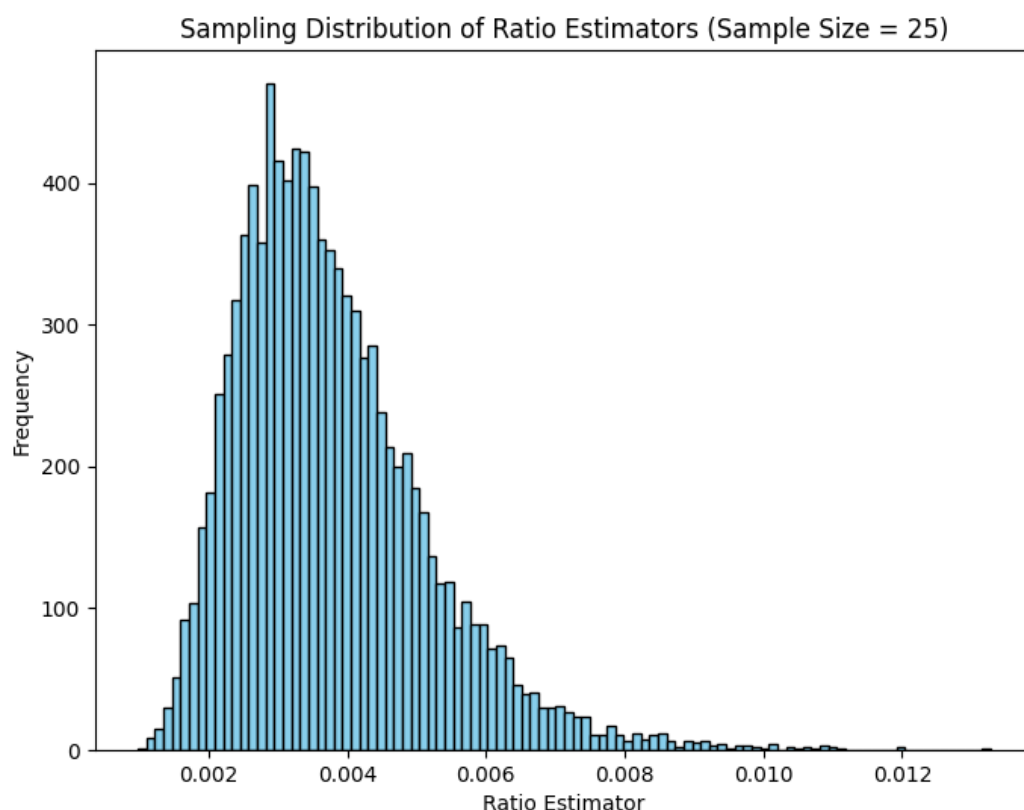


Figure 13: Caption

### 2.1.10 Stratify the counties into four strata and estimate the Population mean and total mortality

For this section, the data is initially divided into four groups based on population size. Subsequently, the mean population is calculated for each group, and the total mean population is estimated by summing and dividing it into four.

Stratum	Mean Population	Total Mortality
1	1463.33	182.18
2	4241.00	280.93
3	8807.83	915.04
4	26702.33	2452.15
Total Mean Population: 10303.62		

### 2.1.11 County Stratification, Sampling Fractions, and Variance Comparison in Population Estimation

For this part, the calculation of Sampling Fractions and Variance is performed through different methods.

## 1. Sampling Fractions:

Calculates the sampling fractions for two allocation methods:

- **Proportional Allocation:** Allocates samples proportionally to each stratum based on population size.
- **Optimal Allocation:** Assumes known variances and allocates samples optimally based on inverse stratum variances.

## 2. Variance Comparison:

Compares the variances of population mean estimates obtained using three sampling methods:

- **Simple Random Sampling:** Randomly selects samples without considering strata.
- **Proportional Allocation:** Allocates samples proportionally to stratum sizes.
- **Optimal Allocation:** Allocates samples optimally based on inverse stratum variances.

Stratum	Proportional Allocation	Optimal Allocation
0	0.043772	$2.323703 \times 10^{-6}$
1	0.098311	$1.003834 \times 10^{-6}$
2	0.205553	$2.243935 \times 10^{-7}$
3	0.652364	$3.508007 \times 10^{-9}$

Table 8: Sampling Fractions for Proportional and Optimal Allocation

Method	Variance
Simple Random Sampling	$1.89888678 \times 10^8$
Proportional Allocation	$1.21515088 \times 10^8$
Optimal Allocation	$3.55543781 \times 10^{-6}$

Table 9: Variances for Different Sampling Methods

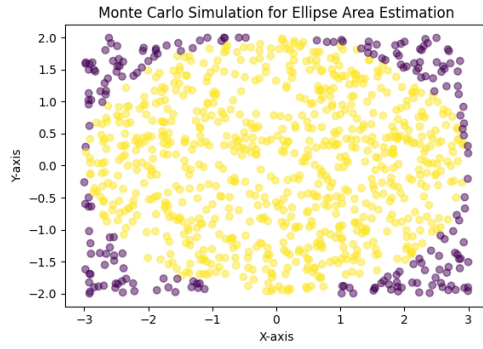
## 2.2 Problem 8

1. This function is written to estimate the value of  $\pi$  using the Monte Carlo method. The Monte Carlo method is used to calculate certain integrals or differential equations that don't have a closed-form solution.

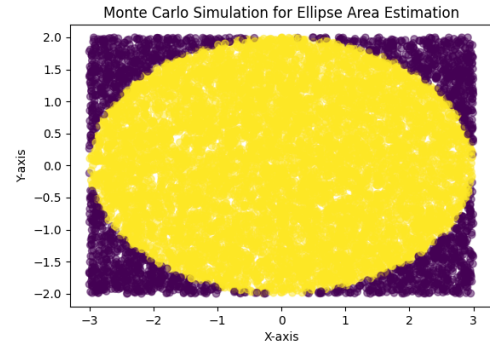
$N$	Estimator Result
10	2.8
100	3.12
1000	3.112
10000	3.148
100000	3.14664
1000000	3.141812
10000000	3.1420716

Table 10: Estimator Results for Different Values of  $N$

2. The decision was made to calculate the area of an ellipse for the Bonus part using the Monte Carlo method. The approach is found to be very similar to a circle estimator, but now involves major and minor diameters. Therefore, a uniform parameter is assigned for the major diameter, and the other is designated for the minor diameter.



(a) Plot for N=100



(b) Plot for N=1000

Figure 14: Ellipse area estimator

## 2.3 Bonus Problem

The mean of Z can be calculated using the law of total expectation:

$$E[Z] = E[Z | X \leq T] \times P(X \leq T) + E[Z | X > T] \times P(X > T) \quad (1)$$

$$E[Z | X \leq T] = E[X | X \leq T] = \alpha \left(1 - e^{-\frac{T}{\alpha}}\right) \quad (2)$$

$$P(X \leq T) = 1 - e^{-\frac{T}{\alpha}} \quad (3)$$

$$E[Z | X > T] = T \quad (4)$$

$$P(X > T) = e^{-\frac{T}{\alpha}} \quad (5)$$

Therefore:

$$E[Z] = \alpha \left(1 - e^{-\frac{T}{\alpha}}\right)^2 + T \times e^{-\frac{T}{\alpha}} \quad (6)$$

The variance of Z can be calculated using the formula:

$$\text{Var}(Z) = E(Z^2) - [E(Z)]^2 \quad (7)$$

$$E(Z^2) = 2\alpha^2 - 3\alpha^2 e^{-\frac{T}{\alpha}} + 2\alpha^2 e^{-\frac{T}{\alpha}} + T^2 e^{-\frac{T}{\alpha}} \quad (8)$$

Being simplified and expanded upon:

$$\text{Var}(Z) = \alpha^2 e^{-\frac{T}{\alpha}} - 2\alpha T e^{-\frac{T}{\alpha}} + T^2 e^{-\frac{T}{\alpha}} \quad (9)$$