

Sampling Distributions

The field of statistical inference consists of those methods used to make decisions or to draw conclusions about a **population**. These methods utilize the information contained in a **sample** from the population in drawing conclusions.

Statistical Inference may be divided into two major areas:

- **Parameter Estimation**
- **Hypothesis testing**

Suppose that we want to obtain a point estimate of a population parameter. We know that before the data is collected, the observations are considered to be **random variables**, say X_1, X_2, \dots, X_n . Therefore, any function of the observation, or any **statistic**, is also a **random variable**.

For example, the sample mean \bar{X} and the sample variance S^2 are **statistics** and they are also **random variables**.

Since a **statistic is a random variable**, it has a probability distribution. We call the probability distribution of a statistic a **sampling distribution**.

In general, if X is a random variable with probability distribution $f(x)$, characterized by the unknown parameter θ , if X_1, X_2, \dots, X_n is a random sample of size n from X the statistic

$$\hat{\Theta} = h(X_1, X_2, \dots, X_n)$$

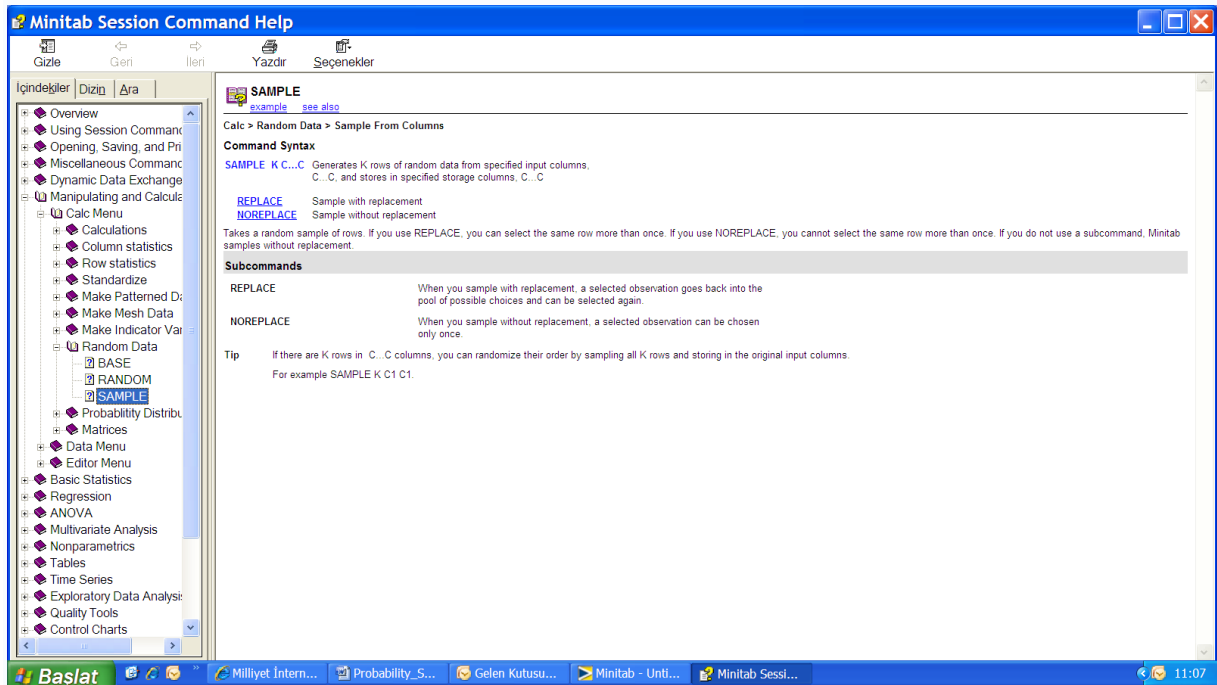
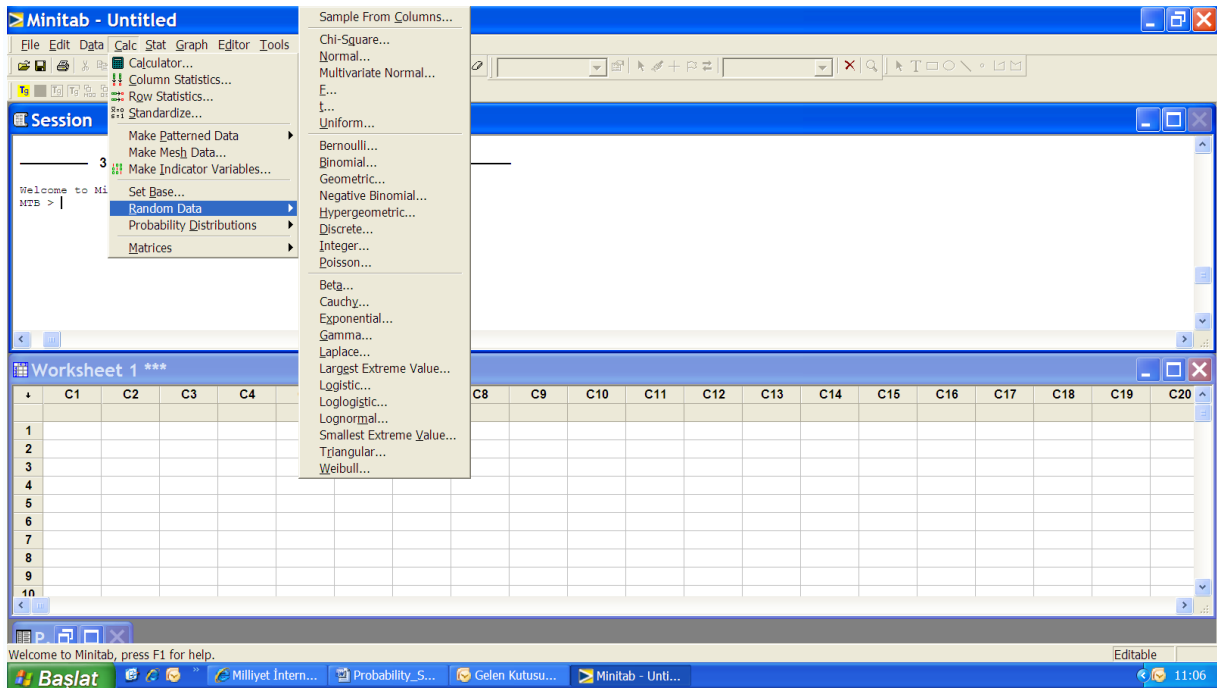
is called a point estimator of θ .

- **Note that** $\hat{\Theta}$ is a random variable because it is a function of random variables

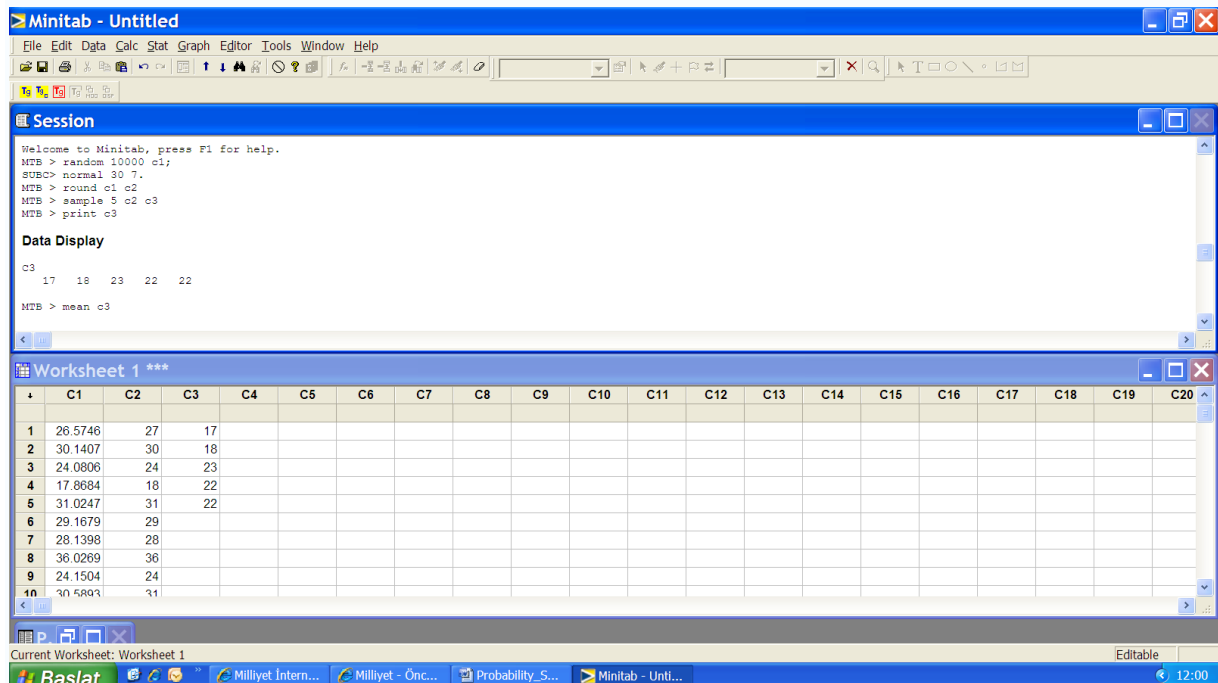
After the sample has been selected, $\hat{\Theta}$ takes on a particular numerical value $\hat{\theta}$ called the point estimate of θ .

Definition: A *point estimate* of some population parameter θ is a single numerical value $\hat{\theta}$ of a statistic $\hat{\Theta}$. The statistic $\hat{\Theta}$ is called the *point estimator*.

Definition: If a sample of n elements is selected from a population of N elements using a sampling plan in which each of the possible samples has the same chance of selection, then the sampling is said to be random and the resulting sample is a *simple random sample*.



As an example, suppose that the random variable X is normally distributed with an unknown mean μ . The sample mean is a point estimator of the unknown population mean μ .



After the sample has been selected, the numerical value \bar{x} is the point estimate of population mean μ .

```

MTB > random 10000 c1;
SUBC> normal 30 7.
MTB > round c1 c2
MTB > sample 5 c2 c3
MTB > print c3

```

Data Display

C3

17 18 23 22 22

MTB > mean c3

Mean of C3

Mean of C3 = 20.4 >>>>>>>>> normal 30 7.

The point estimate of μ is $\bar{x} = \frac{17 + 18 + 23 + 22 + 22}{5} = 20.4$.

In addition to **simple random sampling**, there are other sampling plans that involve randomization and therefore provide probabilistic basis for inference making.

When the population consists of two or more subpopulations, called **strata**, a sampling plan that ensures that each subpopulation is represented in the sample is called a **stratified random sample**.

Definition: *Stratified random sampling involves selecting a simple random sample from each of given number of subpopulations, or **strata**.*

Another form of random sampling is used when the available sampling units are groups of elements, called **clusters**.

Definition: *A cluster sample is a simple random sample of clusters from the available clusters in the population.*

The main difference between **cluster** sampling and **stratified sampling** is that in cluster sampling the cluster is treated as the sampling unit so analysis is done on a population of clusters (at least in the first stage). In stratified sampling, the analysis is done on elements within strata. **In stratified sampling, a random sample is drawn from each of the strata, whereas in cluster sampling only the selected clusters are studied.** The main objective of cluster sampling is to reduce costs by increasing sampling efficiency. This contrasts with stratified sampling where the main objective is to increase precision.

Estimation problems occur frequently in engineering. We often need to estimate

- The mean μ of a single population
- The variance σ^2 (or standard deviation) of a single population
- The proportion p of items in a population that belong to a class of interest
- The difference in means of two populations, $\mu_1 - \mu_2$
- The difference in two population proportions, $p_1 - p_2$.

Reasonable point estimates of these parameters are as follows:

- For μ , the estimate is $\hat{\mu} = \bar{x}$ the sample mean.
- For σ^2 , the estimate is $\hat{\sigma}^2 = S^2$, the sample variance.
- For p , the estimate is $\hat{p} = x/n$ sample proportion, where x is the number of items.
- For $\mu_1 - \mu_2$, the estimate is $\hat{\mu}_1 - \hat{\mu}_2 = \bar{x}_1 - \bar{x}_2$ the difference between the sample means of two independent random samples.
- For $p_1 - p_2$, the estimate is $\hat{p}_1 - \hat{p}_2$, the difference between two sample proportions computed from two independent random samples.

Definition: A *statistic* is any function of the observations in a random sample.

Definition: The probability distribution of a statistic is called a *sampling distribution*.

CENTRAL LIMIT THEOREM

*If random samples of n observations are drawn from a **nonnormal** population with finite mean μ and finite variance σ^2 , then, when n is large, the sampling distribution of the sample mean (sum) \bar{x} is **approximately normally distributed**, with mean μ and standard deviation σ/\sqrt{n} .*

The approximation becomes more accurate as n becomes large **normal 30 7.**

Applet

<http://metalab.uniten.edu.my/~abdrahim/matb344/beaver/diceCLT.html>



```
random 1000 c1-c10 ;
SUBC> integer 1 6.
MTB > rsum c1-c10 c11
MTB > hist c11
```

Histogram of C11

```
MTB > random 100000 c1-c50;
SUBC> integer 1 6.
MTB > rmean c1-c50 c51
MTB > hist c51
```

Histogram of C51

```
MTB > desc c51
```

Descriptive Statistics: C51

Variable	N	N*	Mean	SE Mean	StDev	Minimum	Q1
Median	Q3						
C51	100000	0	3.5006	0.000765	0.2420	2.5200	3.3400
	3.5000	3.6600					

Variable	Maximum
C51	4.6200

Mean and Variance

Suppose X is a discrete uniform random variable on the consecutive integers

$$a, a+1, a+2, \dots, b, \quad \text{for } a \leq b$$

The mean of X is

$$\mu = E(X) = \frac{a+b}{2}$$

The variance of X is

$$\sigma^2 = V(X) = \frac{(b-a+1)^2 - 1}{12}$$

$$\mu = E(X) = \frac{a+b}{2} = \frac{1+6}{2} = 3.5$$

C1	C2							c50	C51
1	3	3	3	1	4	1	3	.	3.00
2	4	4	2	1	3	5	3	.	3.64
3	6	2	3	3	4	5	4	.	3.28
4	2	5	3	6	3	5	6	.	3.88
5	5	6	2	1	4	5	4	.	3.62
4	1	4	2	2	6	2	5	.	3.70
1	5	2	4	5	3	6	3	.	3.74
3	5	6	1	3	3	5	5	.	3.54
2	2	3	6	1	5	5	4	.	3.12
1	2	6	3	2	5	4	6	.	3.36
4	2	2	2	4	3	4	2	.	3.34
3	3	3	4	4	3	2	3	.	3.02
4	1	6	4	5	3	4	2	.	3.44
6	3	2	3	4	5	3	2	.	3.52
4	6	1	5	5	3	6	5	.	3.76
2	4	2	6	6	1	4	4	.	3.04
3	6	2	3	2	1	2	6	.	3.62

$$\sigma^2 = V(X) = \frac{(b-a+1)^2 - 1}{12} = \frac{(6-1+1)^2 - 1}{12} = \frac{35}{12} = 2.916667$$

$$\sigma = \sqrt{2.916667} = 1.7078$$

MTB > stdev c1

Standard Deviation of C1

Standard deviation of C1 = 1.70348

MTB > stdev c2

Standard Deviation of C2

Standard deviation of C2 = 1.70755

.

MTB > stdev c50

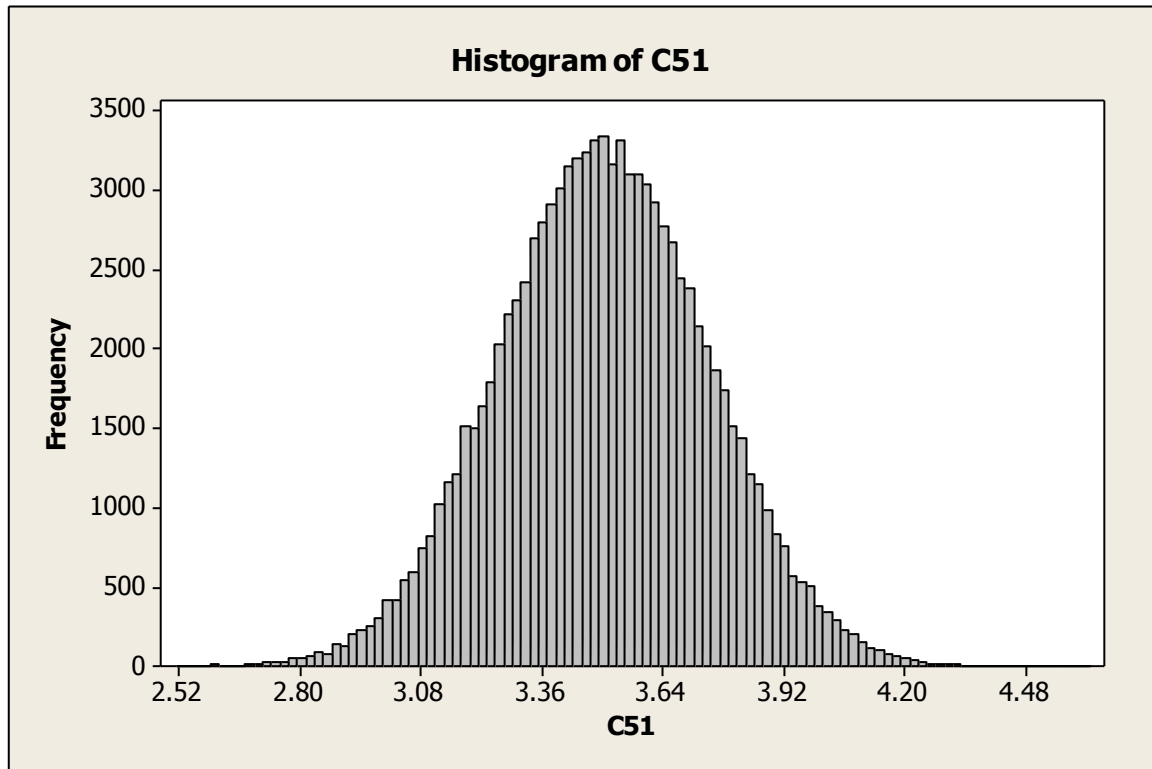
Standard Deviation of C50

Standard deviation of C50 = 1.70769

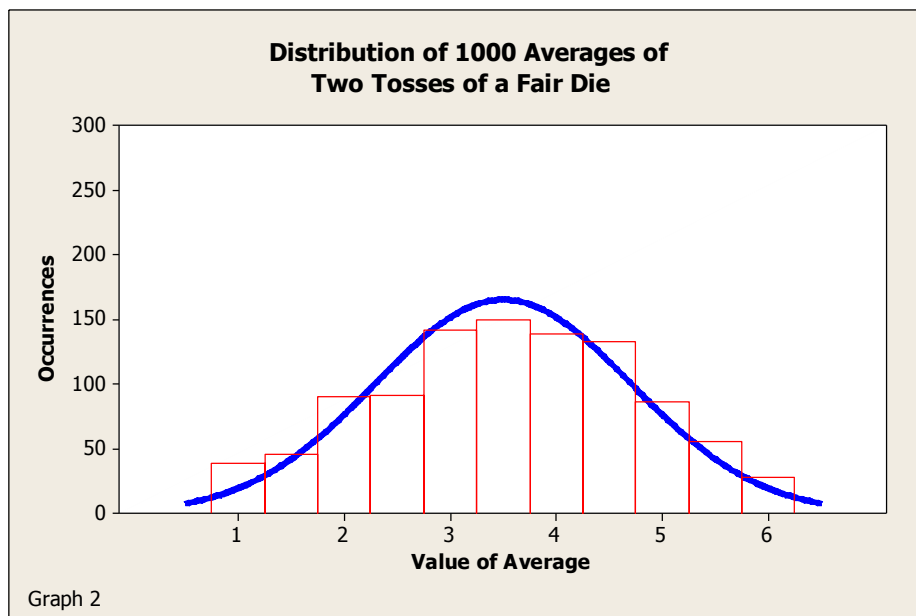
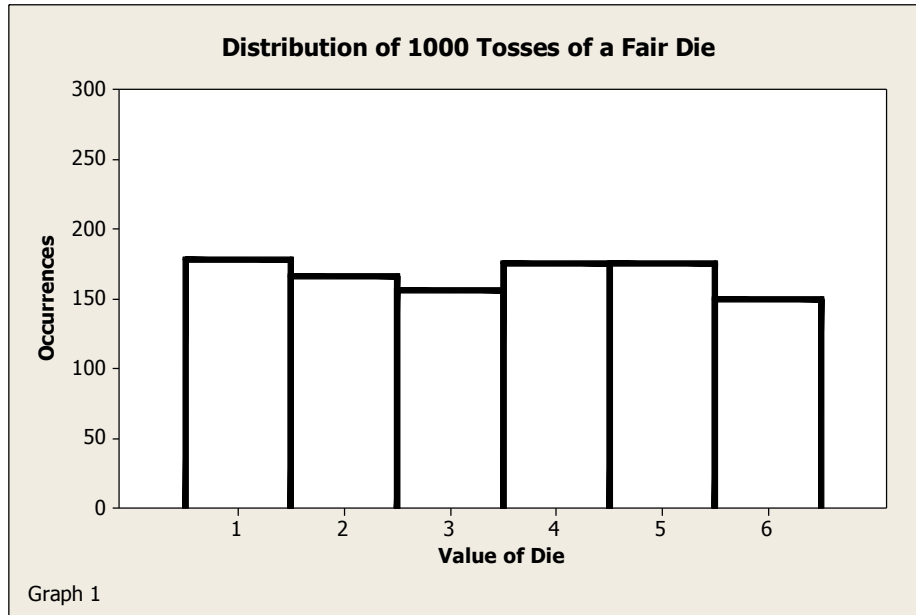
n=50

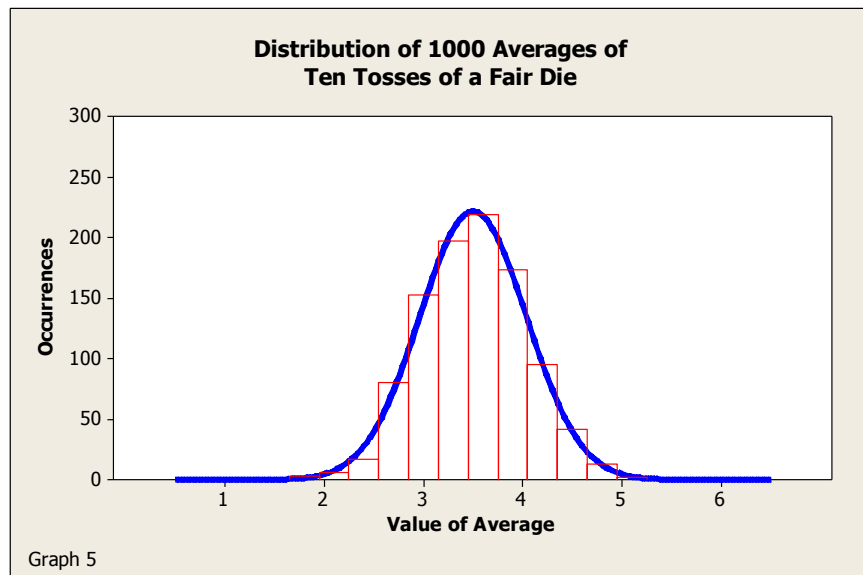
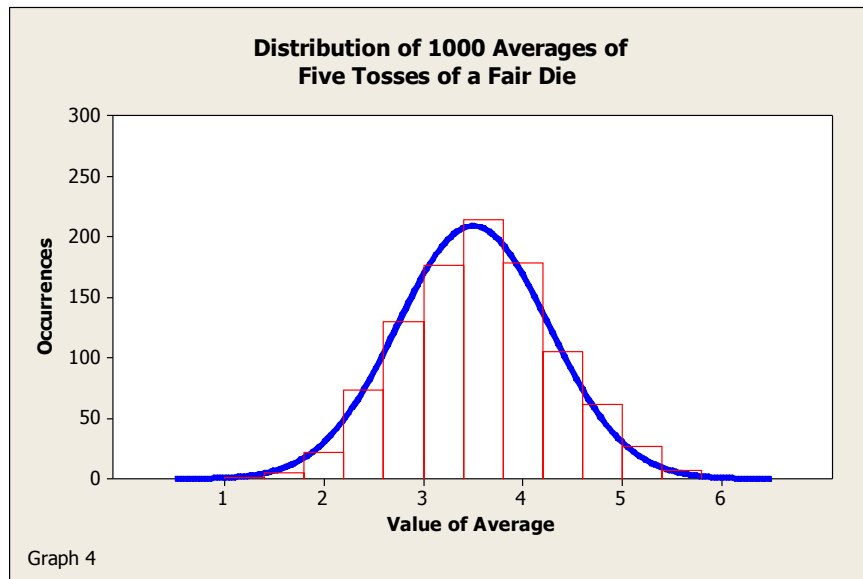
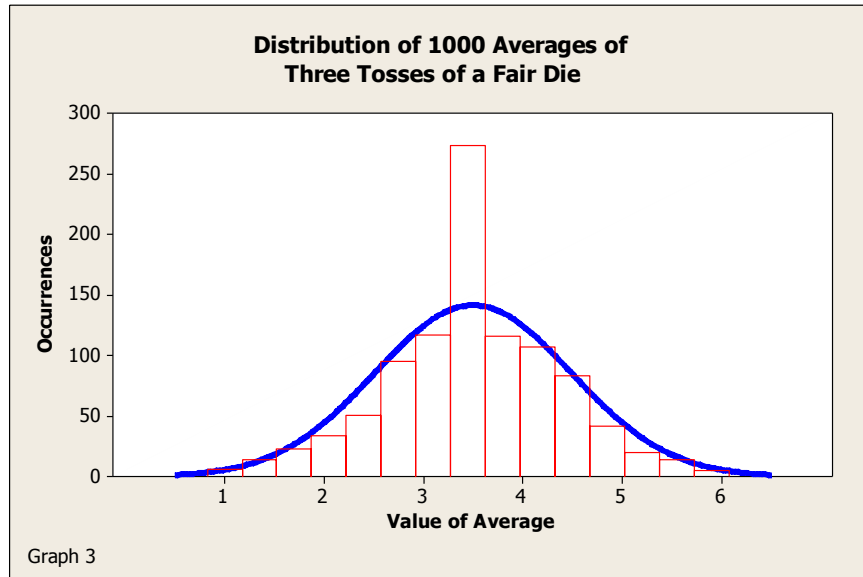
$$\sigma / \sqrt{50} = 1.7078 / \sqrt{50} = 0.24152$$

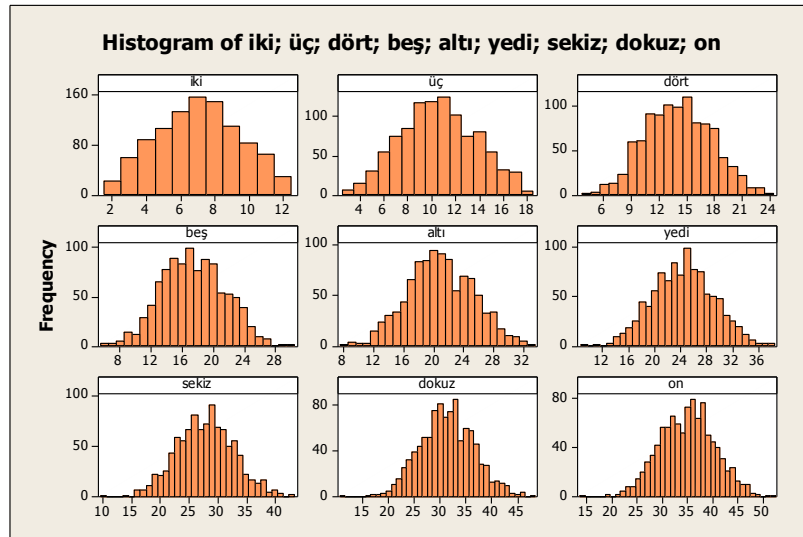
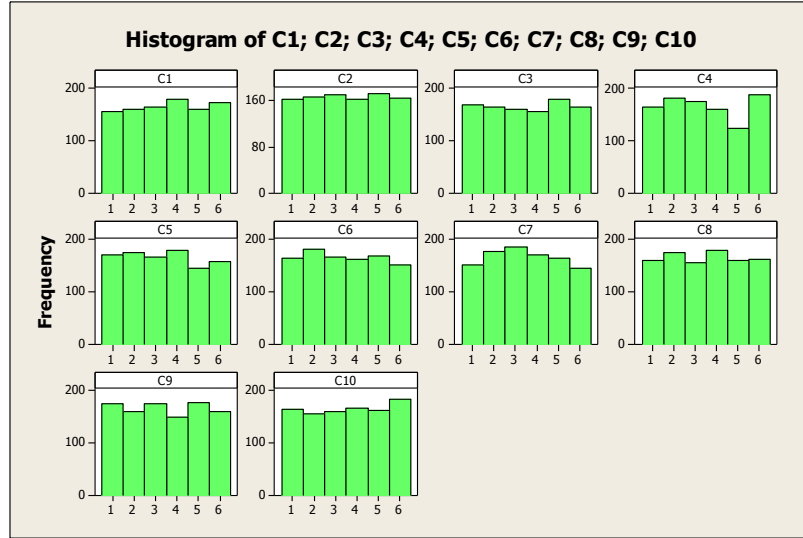
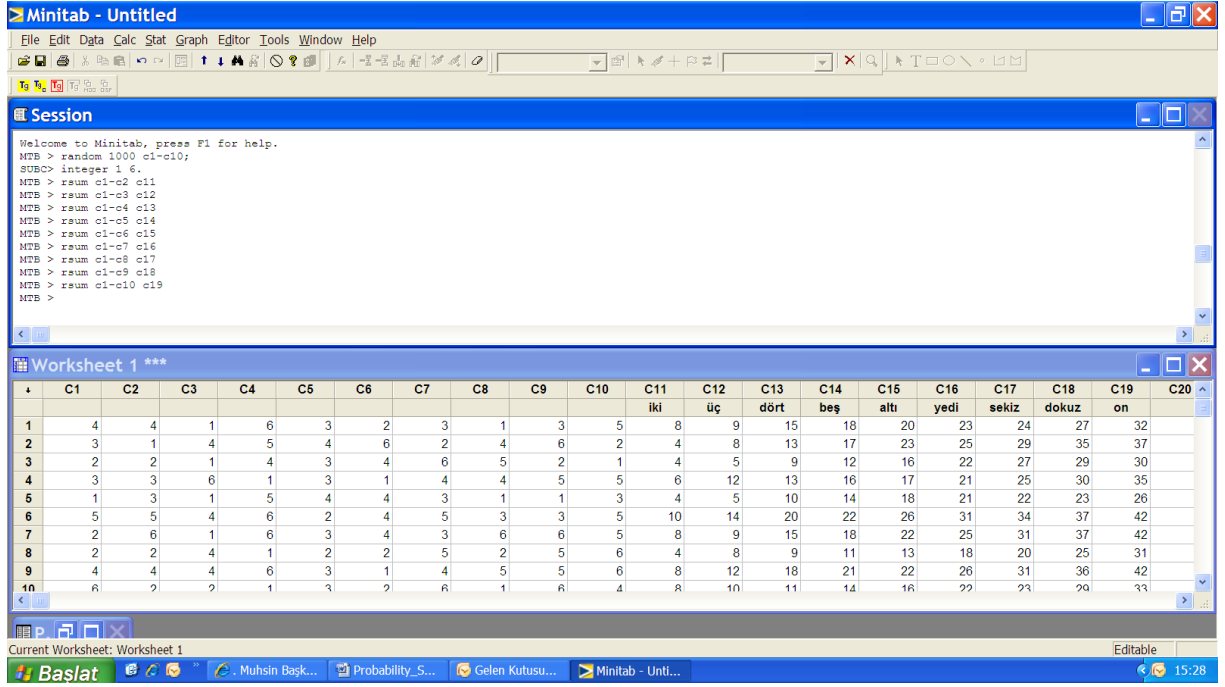
$$\sigma / \sqrt{100000} = 0.24152 / \sqrt{100000} = 0.000765$$

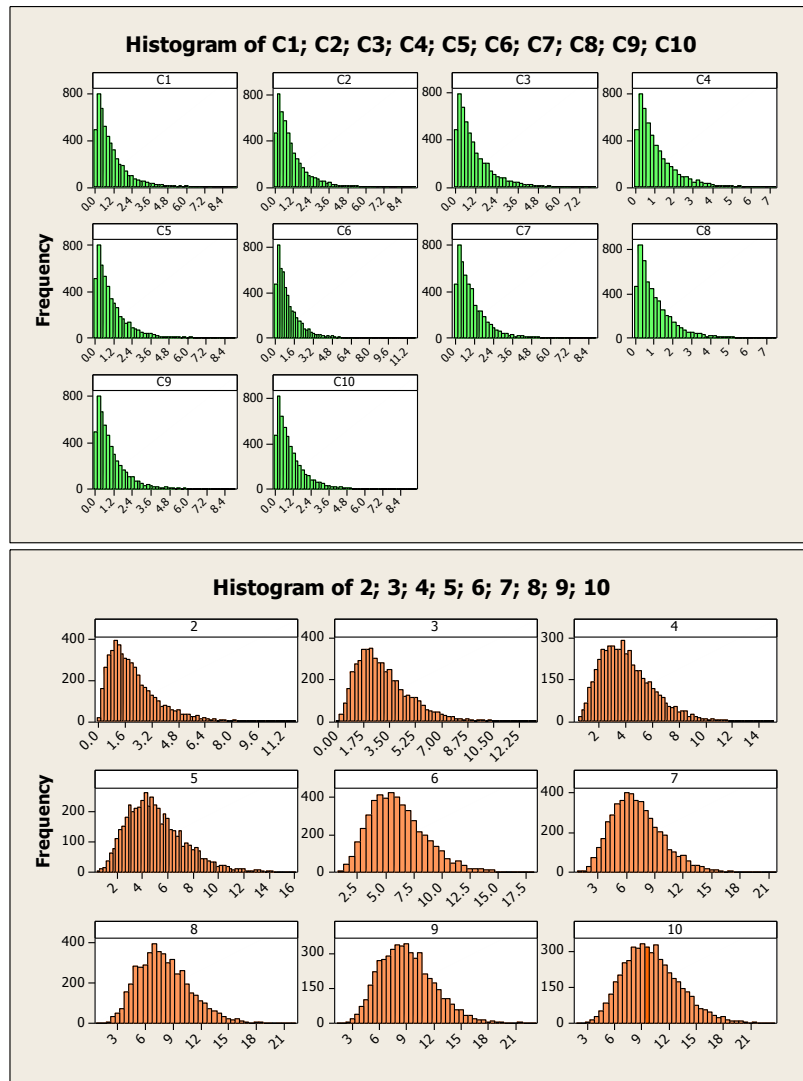
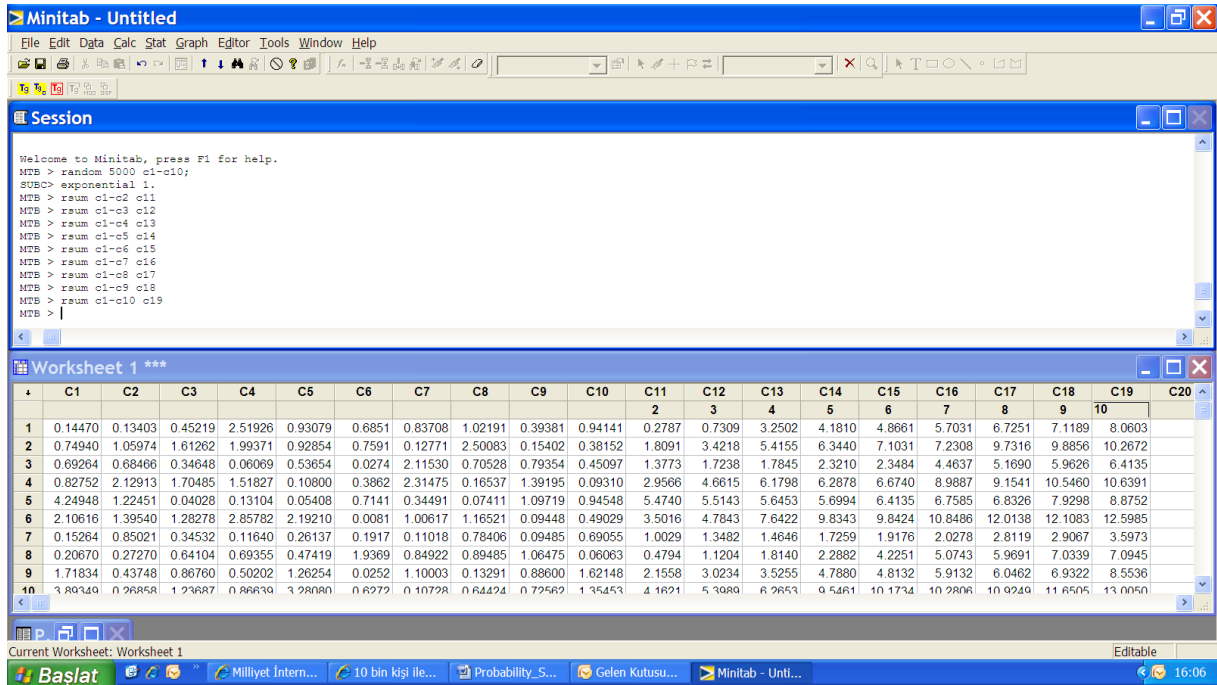


Minitab: Central Limit Theorem

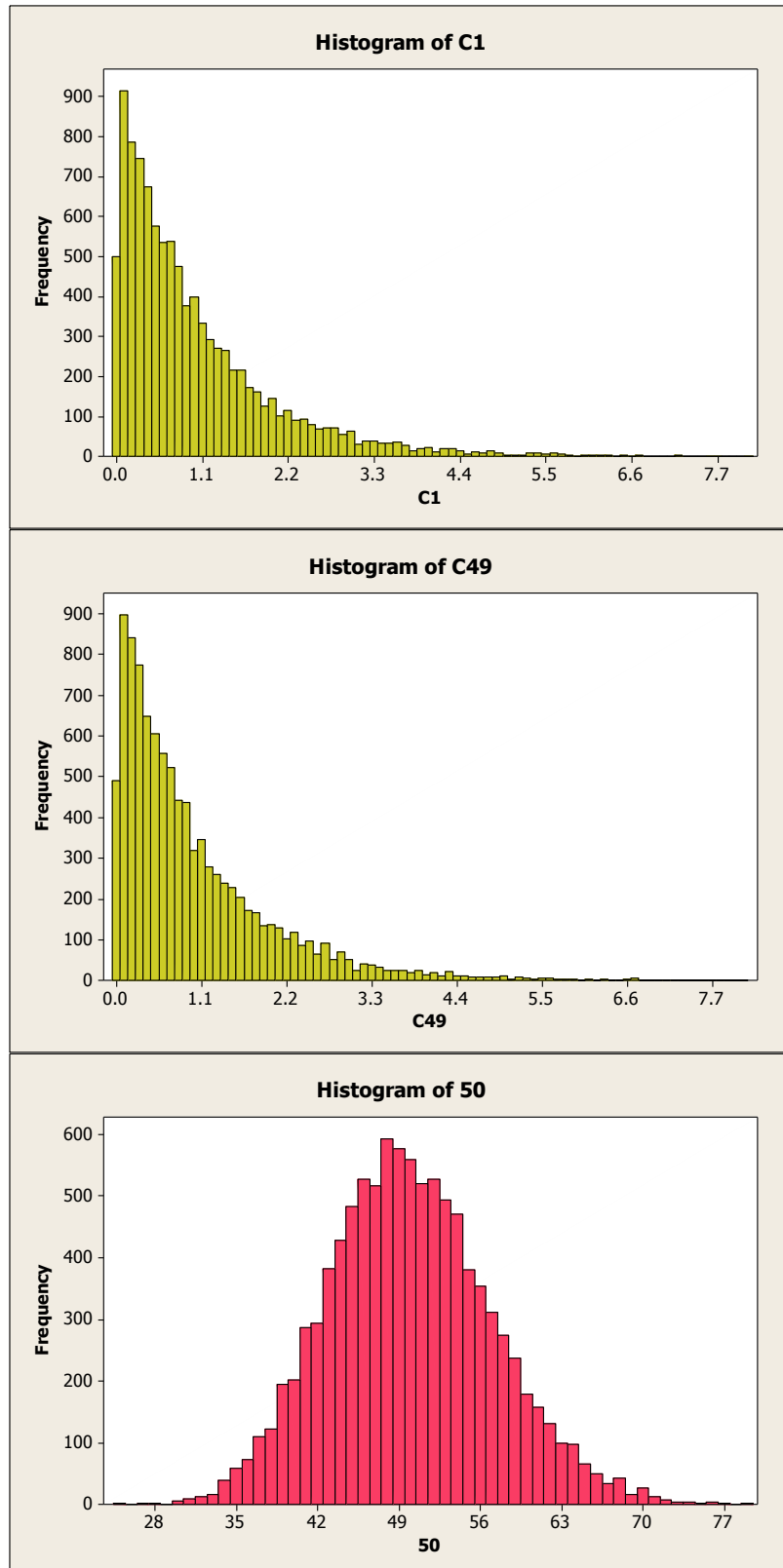




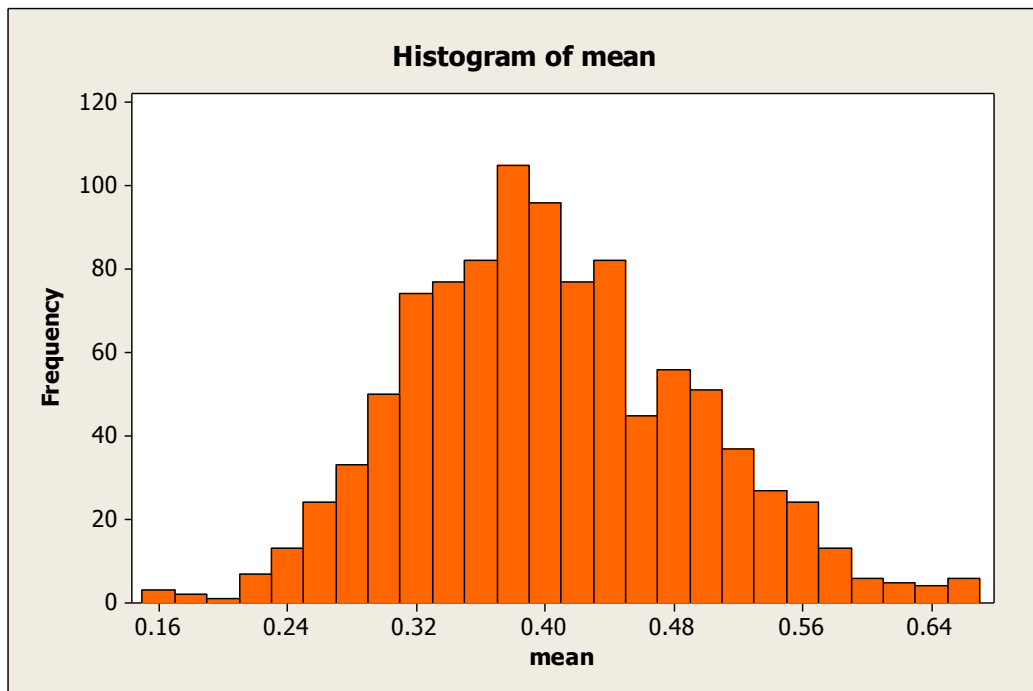
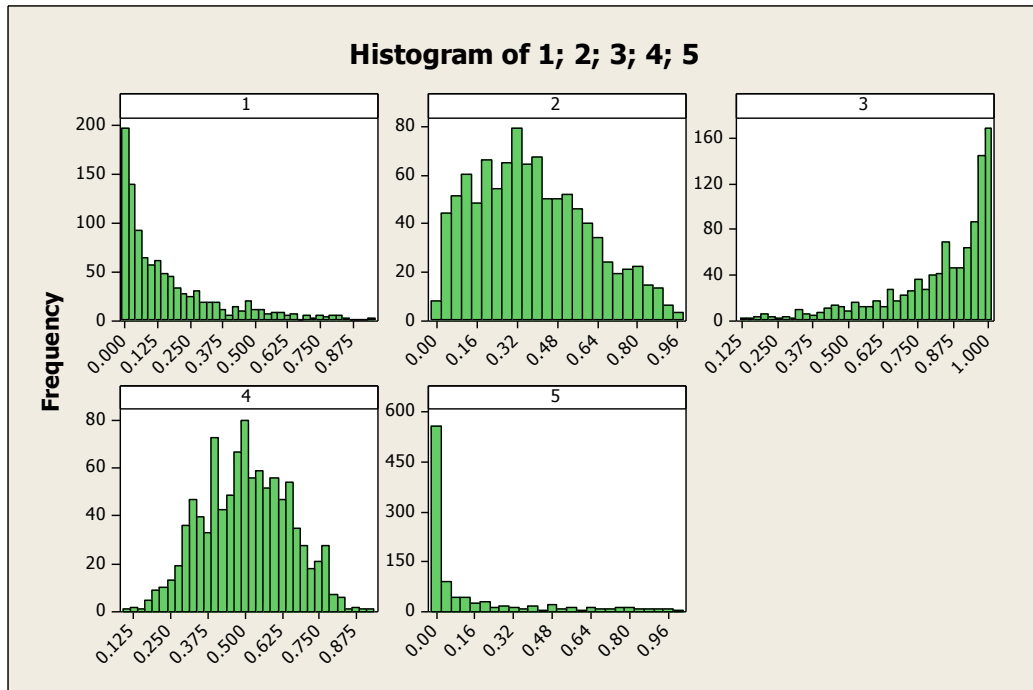




```
MTB > random 10000 c1-c50;
SUBC> expo 1.
MTB > rsum c1-c50 c51
MTB >
```



Different distributions



- If the sampled population is **normal**, then the sampling distribution of \bar{x} will also be normal, no matter what sample size we choose.
- When the sampled population is **approximately symmetric**, the sampling distribution of \bar{x} becomes approximately normal for relatively small values of n. (Dice Example)
- When the sampled population is **skewed**, the sample size n must be larger with n at least 30 before the sampling distribution of \bar{x} becomes approximately normal.

If $n \geq 30$, the normal approximation will satisfactory regardless of the shape of the population.