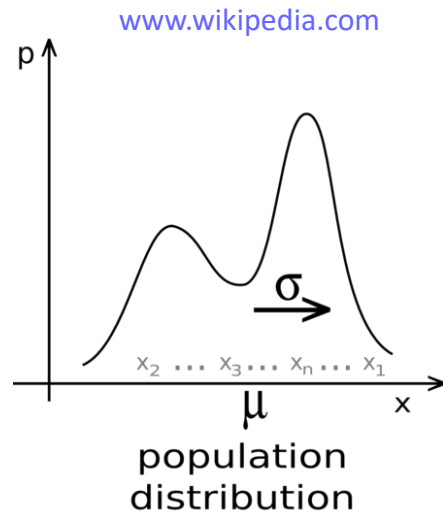


# Central Limit Theorem



“The *distribution of sample means* ( $\bar{y}_1, \bar{y}_2, \bar{y}_3, \bar{y}_4, \dots, \bar{y}_k$ ) follows *a normal distribution*, even when the original variable  $y$  is **NOT** normally distributed.”

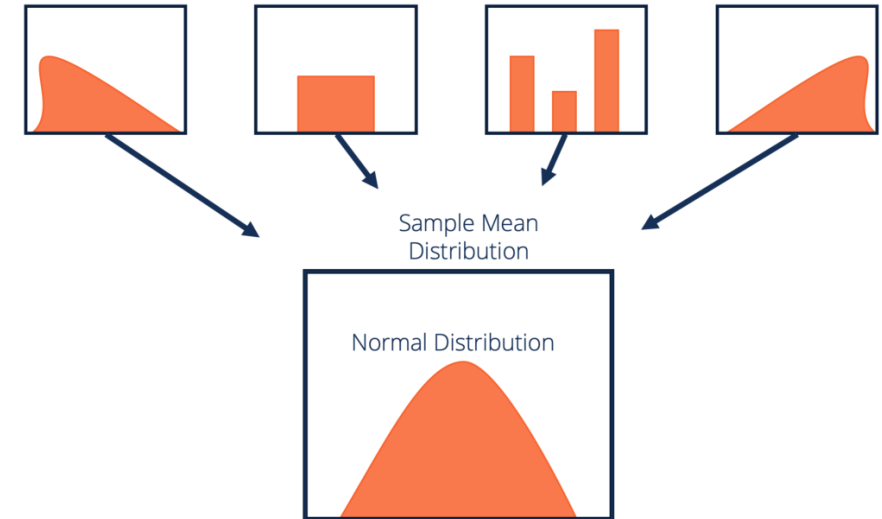
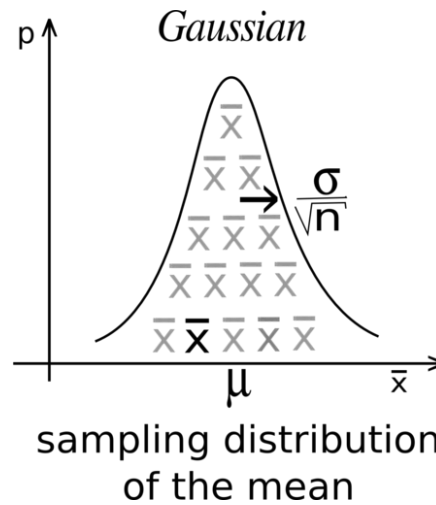
[www.moriah.com](http://www.moriah.com)



samples  
of size  $n$

$\bar{x}$

$\bar{x}$



- What is the mean of distribution of sample means?
- What is the variance of distribution of sample means?

**NOTE:** You do NOT have permission to share this file or any of its contents with anyone else, and/or upload it on internet or any of the platforms where it can be accessed by others.



*“In non-mathematical language, the “CLT” says that whatever the PDF of a variable is, if we randomly sample a “large” number (say  $k$ ) of independent values from that random variable, the sum or mean of those  $k$  values, if collected repeatedly, will have a Normal distribution.*

*It takes some extra thought to understand what is going on here. The process I am describing here takes a sample of (independent) outcomes, e.g., the weights of all of the rats chosen for an experiment, and calculates the mean weight (or sum of weights). Then we consider the less practical process of repeating the whole experiment many, many times (taking a new sample of rats each time). If we would do this, the CLT says that a histogram of all of these mean weights across all of these experiments would show a Gaussian shape, even if the histogram of the individual weights of any one experiment were not following a Gaussian distribution.*

*By the way, the distribution of the means across many experiments is usually called the sampling distribution of the mean.”*

- Seltman, Howard J. "Experimental design and analysis." (2012)

**NOTE:** You do NOT have permission to share this file or any of its contents with anyone else, and/or upload it on internet or any of the platforms where it can be accessed by others.

# Distribution of Sample Means



- The Central Limit Theorem is the explanation for why many real-world random variables tend to have a Gaussian distribution. It is also the justification for assuming that if we could repeat an experiment many times, any sample mean that we calculate once per experiment would follow a Gaussian distribution over the many experiments.
- Since the *distribution of the sample means* with mean ( $\mu$ ) and variance ( $\sigma_y^2/n$ ) follows a normal distribution, then the relationship between the distribution of sample means and the z-distribution is given by:

$$z = \frac{\bar{y} - \mu}{\frac{\sigma_y}{\sqrt{n}}}$$

- What does it tell about value of a random sample mean?
- But, we often don't know the population standard deviation ( $\sigma_y$ ) or variance (!)
- **Can we estimate them?**

**NOTE:** You do NOT have permission to share this file or any of its contents with anyone else, and/or upload it on internet or any of the platforms where it can be accessed by others.

# Estimators of Population



Ref: Montgomery, D. C., "Design and Analysis of Experiments," 8th Ed.,

We may easily show that  $\bar{y}$  and  $S^2$  are unbiased estimators of  $\mu$  and  $\sigma^2$ , respectively.

First consider  $\bar{y}$ . Using the properties of expectation, we have

$$\begin{aligned} E(\bar{y}) &= E\left(\frac{\sum_{i=1}^n y_i}{n}\right) \\ &= \frac{1}{n} \sum_{i=1}^n E(y_i) \\ &= \frac{1}{n} \sum_{i=1}^n \mu \\ &= \mu \end{aligned}$$

$$\begin{aligned} E(S^2) &= E\left[\frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n-1}\right] \\ &= \frac{1}{n-1} E\left[\sum_{i=1}^n (y_i - \bar{y})^2\right] \\ &= \frac{1}{n-1} E(SS) \end{aligned}$$

$$\begin{aligned} E(SS) &= E\left[\sum_{i=1}^n (y_i - \bar{y})^2\right] \\ &= E\left[\sum_{i=1}^n y_i^2 - n\bar{y}^2\right] \\ &= \sum_{i=1}^n (\mu^2 + \sigma^2) - n(\mu^2 + \sigma^2/n) \\ &= (n-1)\sigma^2 \end{aligned}$$

$$E(S^2) = \frac{1}{n-1} E(SS) = \sigma^2$$

$S^2$  is an unbiased estimator of  $\sigma^2$ .

where  $SS = \sum_{i=1}^n (y_i - \bar{y})^2$  is the **corrected sum of squares** of the observations  $y_i$ .

**NOTE:** You do NOT have permission to share this file or any of its contents with anyone else, and/or upload it on internet or any of the platforms where it can be accessed by others.

# Degrees of Freedom (DOF)



If  $y$  is a random variable with variance  $\sigma^2$ ,

and sum of squares  $SS = \sum (y_i - \bar{y})^2$  has 'v' degrees of freedom, then  $E\left(\frac{SS}{v}\right) = \sigma^2$

The number of **degrees of freedom of a sum of squares** is equal to **the number of independent elements** in that sum of squares.

For example,  $SS = \sum (y_i - \bar{y})^2$  is a sum of squares of 'n' elements, i.e.,  $y_1 - \bar{y}, y_2 - \bar{y}, \dots, y_n - \bar{y}$

Note that these 'n' elements are not all independent because  $\sum (y_i - \bar{y}) = 0$

Therefore, only n-1 of them are independent, implying that SS has (n-1) degrees of freedom.

$$E\left(\frac{SS}{n-1}\right) = \sigma^2$$

**NOTE:** You do NOT have permission to share this file or any of its contents with anyone else, and/or upload it on internet or any of the platforms where it can be accessed by others.

# Consequence of CLT



If  $y_1, y_2, y_3 \dots, y_n$  is a sequence of 'n' independent and identically distributed random variables with  $E(y_i) = \mu$  and  $V(y_i) = \sigma^2$  (both finite)

If we define,  $x = y_1 + y_2 + y_3 + \dots + y_n$  Then what is the distribution of 'x' as 'n' becomes sufficiently large?  
 $= n \bar{y}$

In other words,  $Z_n = \frac{x - n\mu}{\sqrt{n\sigma^2}}$  is a standard normal distribution as  $n \rightarrow \infty$

IMP!

**'sum of n independent and identically distributed random variables is approximately normally distributed'**

Frequently, we think of the error in an experiment as arising in an additive manner from several independent sources; consequently, **the normal distribution becomes a reasonable model for the combined experimental error.**

**NOTE:** You do NOT have permission to share this file or any of its contents with anyone else, and/or upload it on internet or any of the platforms where it can be accessed by others.

# z-Distribution (Std. Normal PDF)

If  $y_1, y_2, \dots, y_n$  is a random sample from the **ANY** distribution, then

$$z = \frac{\bar{y} - \mu}{\sigma / \sqrt{n}}$$

is distributed as **Standard Normal Distribution, i.e., NPDF (0, 1)**

## Standard Normal Probabilities

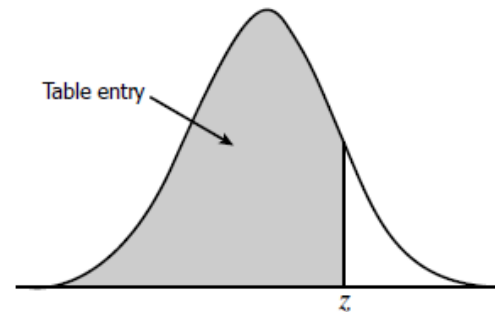


Table entry for  $z$  is the area under the standard normal curve to the left of  $z$ .

$z$	.00	.01	.02	.03	.04	.05	.06	.07	.08	.09
0.0	.5000	.5040	.5080	.5120	.5160	.5199	.5239	.5279	.5319	.5359
0.1	.5398	.5438	.5478	.5517	.5557	.5596	.5636	.5675	.5714	.5753
0.2	.5793	.5832	.5871	.5910	.5948	.5987	.6026	.6064	.6103	.6141
0.3	.6179	.6217	.6255	.6293	.6331	.6368	.6406	.6443	.6480	.6517
0.4	.6554	.6591	.6628	.6664	.6700	.6736	.6772	.6808	.6844	.6879

**DIY** Read Chapter 2 Pages 32- 52 from Design and Analysis of Experiments, 8<sup>th</sup> Ed.

**NOTE:** You do NOT have permission to share this file or any of its contents with anyone else, and/or upload it on internet or any of the platforms where it can be accessed by others.

# t-Distribution



If  $y_1, y_2, \dots, y_n$  is a random sample from the  $N(\mu, \sigma^2)$  distribution, then

$$t = \frac{\bar{y} - \mu}{S/\sqrt{n}}$$

is distributed as  $t$  with  $n - 1$  degrees of freedom.

## Side Notes:

What happens when  $n \rightarrow \infty$ ?

$$\Gamma(z) = \int_0^\infty t^{z-1} e^{-t} dt, \quad \Re(z) > 0$$

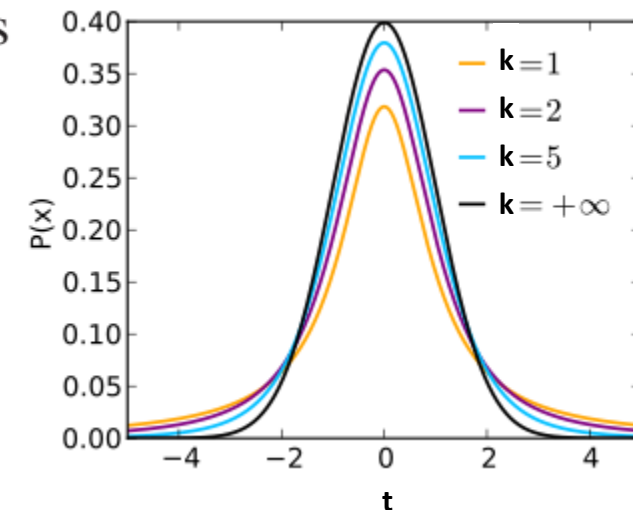
**$t$  distribution with  $k$  degrees of freedom**, denoted  $t_k$ . The density function of  $t$  is

$$f(t) = \frac{\Gamma[(k+1)/2]}{\sqrt{k\pi}\Gamma(k/2)} \frac{1}{[(t^2/k) + 1]^{(k+1)/2}} \quad -\infty < t < \infty$$

**Note:** The increased spread reflects the added uncertainty *due to unknown  $\sigma_y$*  that gets estimated by  $s$ , which itself is prone to sampling errors. **What will happen as  $k$  increases?**

**DIY**

Read Chapter 2 Pages 32- 52 from Design and Analysis of Experiments, 8<sup>th</sup> Ed.



**NOTE:** You do NOT have permission to share this file or any of its contents with anyone else, and/or upload it on internet or any of the platforms where it can be accessed by others.



# Chi-Square ( $\chi^2$ ) Distribution



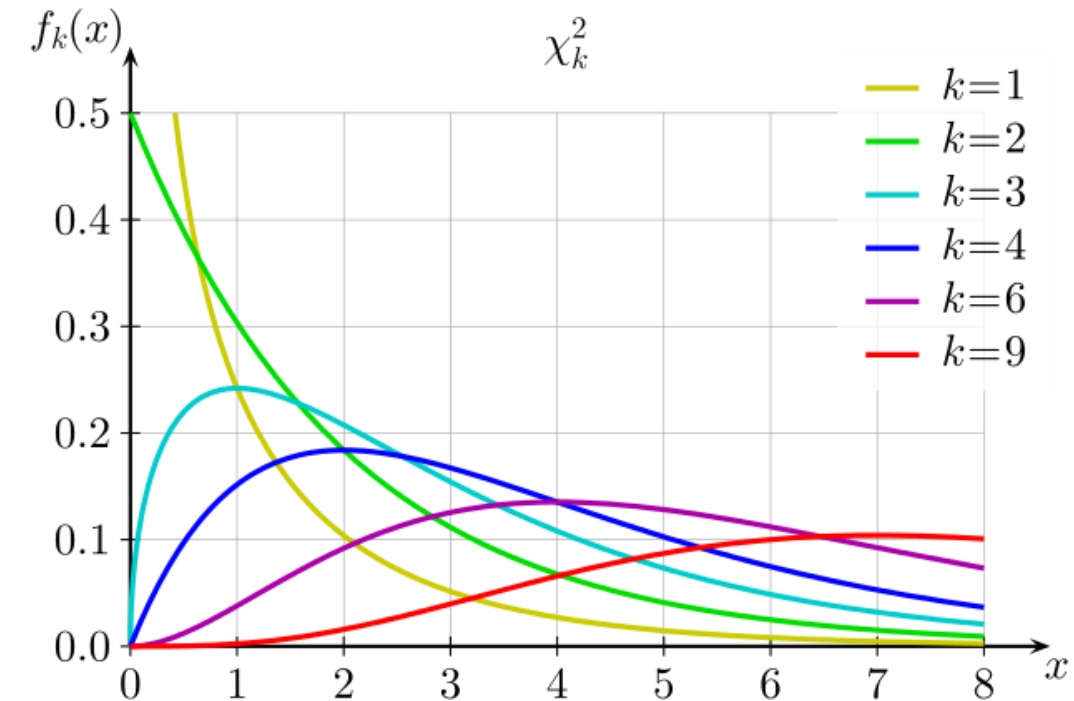
If  $z_1, z_2, z_3 \dots, z_k$  are **normally and independently distributed** random variables with mean 0 and variance 1 [NID (0,1)]

And if we define,  $\chi = z_1^2 + z_2^2 + \dots + z_k^2$

Then 'x' follows the **chi-square distribution with k degrees of freedom**

$$f(x; k) = \begin{cases} \frac{x^{\frac{k}{2}-1} e^{-\frac{x}{2}}}{2^{\frac{k}{2}} \Gamma\left(\frac{k}{2}\right)}, & x > 0; \\ 0, & \text{otherwise.} \end{cases}$$

- The distribution is asymmetric or skewed
- Mean,  $\mu = k$  and Variance,  $\sigma^2 = 2k$



**NOTE:** You do NOT have permission to share this file or any of its contents with anyone else, and/or upload it on internet or any of the platforms where it can be accessed by others.

# Chi-Square ( $\chi^2$ ) Distribution



- What would be an example of chi-square distribution?

- Remember

$SS = \sum_{i=1}^n (y_i - \bar{y})^2$  is the corrected sum of squares of the observations  $y_i$ .

$$DOF = n - 1$$

$$\underline{E(S^2)} = \frac{1}{n-1} E(SS) = \underline{\sigma^2}$$

and we see that  $S^2$  is an unbiased estimator of  $\sigma^2$ .

$$\boxed{\sum y_i - \bar{y} = 0}$$

$$\sum \frac{(y_i - \bar{y})^2}{\sigma^2}$$

$\left( \frac{y_1 - \bar{y}}{\sigma} \right)^2, \frac{y_2 - \bar{y}}{\sigma} + \dots + \frac{y_n - \bar{y}}{\sigma}$

$$\frac{SS}{\sigma^2} = \frac{\sum_{i=1}^n (y_i - \bar{y})^2}{\sigma^2} \sim \chi_{n-1}^2$$

$$SS = \sigma^2 \chi_{n-1}^2$$

Sample variance,  $S^2 = \frac{SS}{n-1}$

Therefore, if the observations in the sample are NID  $(\mu, \sigma^2)$ , then the distribution of  $S^2$  is  $\left( \frac{\sigma^2}{n-1} \right) \chi_{n-1}^2$

Thus, the sampling distribution of the sample variance is a constant times the chi-square distribution if the population is normally distributed.

**NOTE:** You do NOT have permission to share this file or any of its contents with anyone else, and/or upload it on internet or any of the platforms where it can be accessed by others.

# F-Distribution



If  $\chi_u^2$  and  $\chi_v^2$  are two independent chi-square random variables with  $u$  and  $v$

degrees of freedom, respectively, then, the ratio  $F_{u,v} = \frac{\chi_u^2/u}{\chi_v^2/v}$

follows a F-distribution with  $u$  degrees of freedom of numerator and  $v$  degrees of freedom of denominator

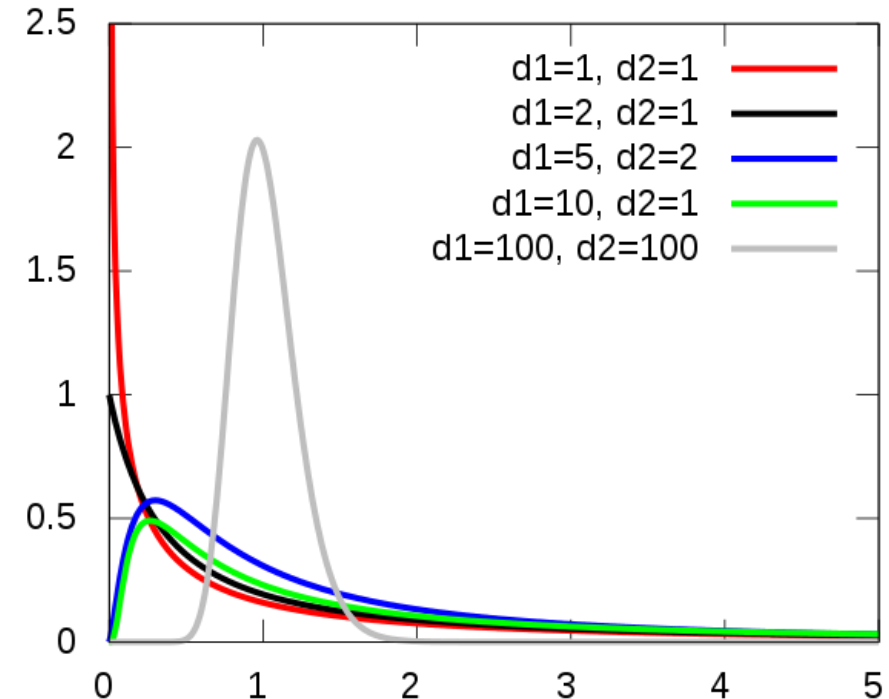
$$F(x) = \frac{\Gamma\left(\frac{u+v}{2}\right)\left(\frac{u}{v}\right)^{u/2} x^{(u/2)-1}}{\Gamma\left(\frac{u}{2}\right)\Gamma\left(\frac{v}{2}\right)\left[\left(\frac{u}{v}\right)x + 1\right]^{(u+v)/2}} \quad 0 < x < \infty$$

## Example of F-distribution

Suppose we have *two independent normal populations* with common variance  $\sigma^2$ .

If  $y_{11}, y_{12}, y_{13}, \dots, y_{1n_1}$  is a random sample of  $n_1$  observations from the first population and  $y_{21}, y_{22}, y_{23}, \dots, y_{2n_2}$  is a

random sample of  $n_2$  observations from the second population, **Then,**  $\frac{S_1^2}{S_2^2} \sim F_{n_1-1, n_2-1}$



**NOTE:** You do NOT have permission to share this file or any of its contents with anyone else, and/or upload it on internet or any of the platforms where it can be accessed by others.