# Central limit Theorem and Confidence Intervals

Sample from population

$$X_1, X_2, \ldots, X_n \quad \text{i.i.d.} \quad X$$

[ with replacement ; without replacement ]

- Suppose $\mu = E[x]$ , $\sigma^2 = \text{Var}[x]$

$$\bar{X} = \frac{X_1 + X_2 + \cdots + X_n}{n} \equiv \cdots \text{estimate} \cdots = \mu$$

[Justification] required

$$\sigma_x^2 = \sum_{i=1}^{n} \frac{(X_i - \bar{X})^2}{n-1} \equiv \cdots \text{estimate} \cdots = \sigma^2$$

[Justification] required

- - Summary : $\bar{X}, \sigma_x^2$ , median, min, max
    of the distribution
  - Plots : Histogram, box plots, qr-qr plots.

- Empirical distribution :- $S = \{X_1, X_2, \ldots, X_n\}$
    - include repeat observations

$$f_n (t) = \frac{\#\{i : X_i = t\}}{n}$$

  - Discrete p.m.f. on $S$ $\equiv$ inference based on
    this is called descriptive statistics

    Ex: Suppose $Y$ (r.v) has p.m.f $f_n(\cdot)$
    $$P(Y=t) = f_n(t)$$
    $$E[Y] = \bar{X} \quad , \quad \text{Var}[Y] = ?$$

$Z \sim \text{Normal} (0,1)$ if

$$P(Z \leq z) = \int_{-\infty}^{z} \frac{e^{-\frac{x^2}{2}}}{\sqrt{2\pi}} dx$$

Normal Tables to evaluate numerically

$Z$ has p.d.f
$$f_z(z) = \frac{e^{-\frac{z^2}{2}}}{\sqrt{2\pi}} \quad z \in \mathbb{R}$$

Central limit Theorem
- Distribution occurs naturally.

- arises as sum of independent processes

- # of leafs in a tree
- height of individuals.
[Suitable interpretation]

# Normal Distribution: PDF

You can calculate the values of the normal density function using the the `dnorm` command.

```
> dnorm(0)
[1] 0.3989423
> dnorm(1)
[1] 0.2419707
> dnorm(0, mean=4, sd=3)
[1] 0.05467002
```
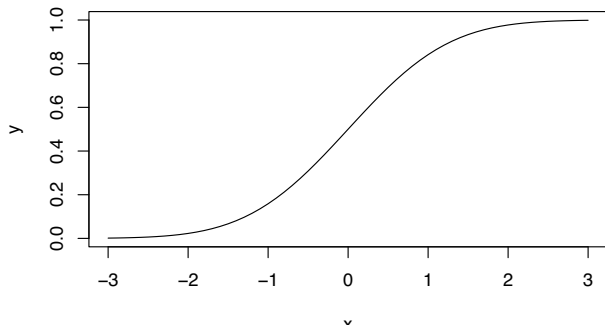
# Normal Distribution: CDF

You can calculate the values of the cummulative distribution function of the normal using the the pnorm command.
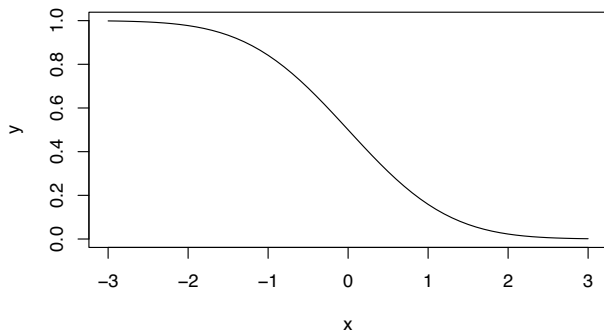
```
> pnorm(0)
> pnorm(1)
> x = seq(-3,3, by=0.1); y = pnorm(x) ;plot(x,y, type="l")
```

# Normal Distribution: Tail Probabilities

```
> pnorm(0, lower.tail=FALSE)
> pnorm(1, lower.tail=FALSE)
> x = seq(-3,3, by=0.1); y = pnorm(x, lower.tail=FALSE)
> plot(x,y, type="l")
```
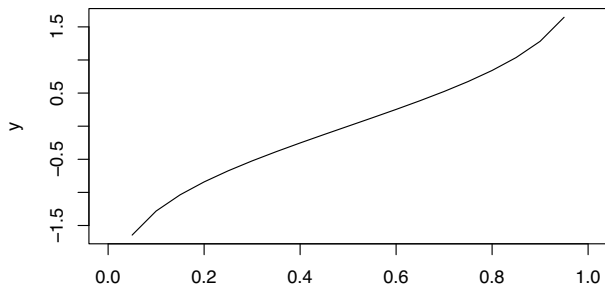
# Normal Distribution: quantiles

```
> qnorm(0.68); qnorm(0.95);qnorm(0.997)
[1] 0.4676988
[1] 1.644854
[1] 2.747781
> x = seq(0,1, by=0.05); y = qnorm(x);plot(x,y, type="l")
```
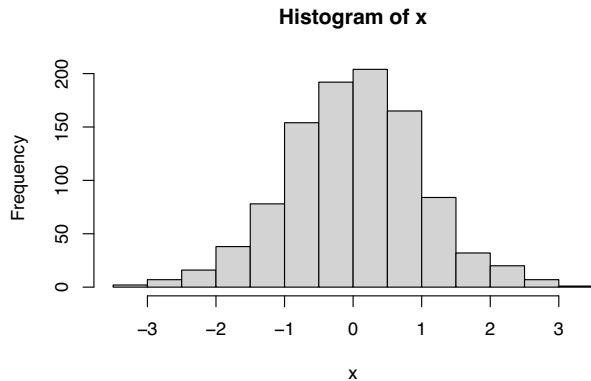
# Normal Distribution: samples

```
> x=rnorm(1000)
> hist(x)
```

**Histogram of x**



— Symmetric
  about
    mean = 0

```
> pnorm(1) - pnorm(-1) # within one standard deviation
[1] 0.6826895
> pnorm(2) - pnorm(-2) # within two standard deviation
[1] 0.9544997
> pnorm(3) - pnorm(-3) # within three standard deviation
[1] 0.9973002
```

# Central Limit Theorem

$$S_n = \sum_{i=1}^{n} X_i \quad \text{with} \quad X_i \sim \text{Bernoulli}(p) \quad \text{independent for } i \geq 1$$

Q:– How good is the Normal approximation?

Suppose each $X_i$ was distributed as Bernoulli $(p)$ random variable. Then $S_n$ is a Binomial$(n,p)$ random variable. Let us check for what $p$ does

$$\frac{S_n - np}{\sqrt{np(1-p)}}$$

is close to a Normal distribution.

Noted in Work sheet $\therefore$ $\left| \, \mathbb{P}\left( \dfrac{S_n - np}{\sqrt{np(1-p)}} \leq x \right) - \mathbb{P}(Z \leq x) \, \right| \leq \dfrac{C_1}{\sqrt{n}}$

$$\longrightarrow 0 \quad \text{as } n \to \infty$$
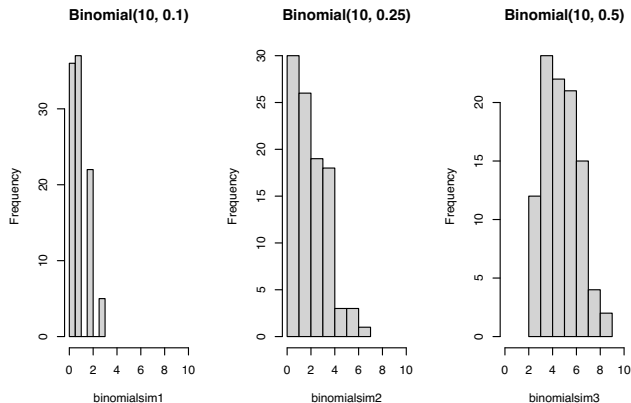
$Z \sim \text{Normal}(0,1)$

# Central Limit Theorem

We may simulate Binomial samples either direclty by `rbinom` command or usi ng the `replicate` and `rbinom` command.

```
> binomialsim1 = rbinom(100,10,0.1)
> # generates 100 Binomial (10,0.1) samples
>
> binomialsim2 = replicate(100, rbinom(1,10,0.25))
> # generates 100 Binomial (10,0.25) samples
>
> binomialsim3 = replicate(100, rbinom(1,10,0.5))
> # generates 100 Binomial (10,0.5) samples
>
```
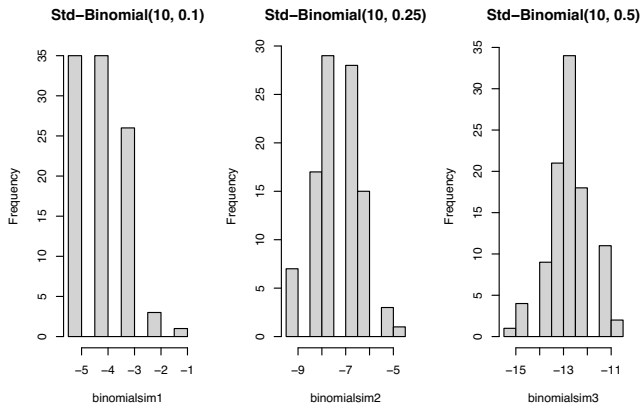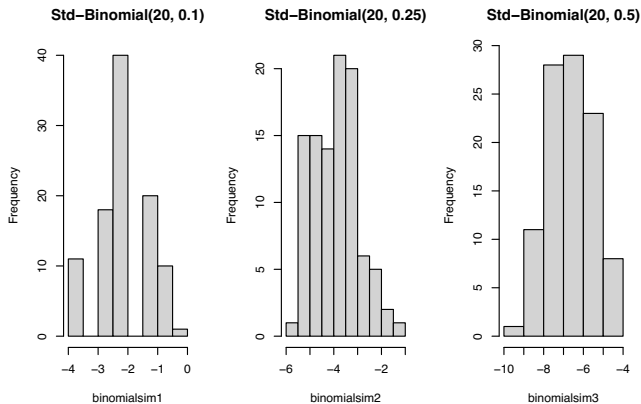
**Binomial(10, 0.1)**   **Binomial(10, 0.25)**   **Binomial(10, 0.5)**

From the above it seems that at $n = 10$ the symmetry is achieved when $p = 0.5$ and not at $p = 0.1$ and $p = 0.25$

**Std−Binomial(10, 0.1)**   **Std−Binomial(10, 0.25)**   **Std−Binomial(10, 0.5)**

Perhaps $n = 10$ is not large enough to see the Central Limit Theorem occuring.

**Std−Binomial(20, 0.1)** **Std−Binomial(20, 0.25)** **Std−Binomial(20, 0.5)**
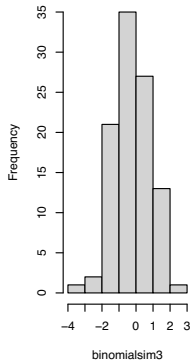
$n = 20$ is better.

**Std−Binomial(50, 0.1)**  **Std−Binomial(50, 0.25)**  **Std−Binomial(50, 0.5)**

$n = 50$ we get closer to Normal distribution

# Role of $n$ versus $p$

Binomial Random variable is close to Normal when the distribution is symmetric. That is when $p$ is close to 0.5. Otherwise the general rule that we can apply is that when

$$np \geq 5 \text{ and } n(1 - p) \geq 5.$$

then Binomial(n,p) is close to Normal distribution.

# Central Limit Theorem — "True in general for Sums"

**Sample:** $X_1, X_2, \ldots, X_n$ i.i.d $X$     $E[x] = \mu$ ; $Var[x] = \sigma^2$

We could rephrase the result as:

> Fundamental Result

Let $X_1, X_2, \ldots$ be i.i.d. random variables with finite mean $\mu$, finite variance $\sigma^2$. Then

$$\frac{(S_n - n\mu)}{\sqrt{n}\sigma} \xrightarrow{d} Z, \tag{3}$$

where $S_n = X_1 + X_2 + \ldots + X_n$ and $Z \sim \text{Normal}(0, 1)$.

i.e. $\left| \mathbb{P}\left( \frac{S_n - n\mu}{\sqrt{n}\,\sigma} \leq x \right) - \mathbb{P}(Z \leq x) \right| \longrightarrow 0$   as $n \to \infty$

"occurs naturally as sums" $\iff$ "$S_n \sim N(n\mu, n\sigma^2)$"

# Central Limit Theorem

$$\frac{S_n - n\mu}{\sqrt{n}\,\sigma} = \frac{\sqrt{n}}{n}\left(\frac{S_n}{n} - \mu\right)\Big/ \sqrt{n}\,\sigma = \sqrt{n}\left(\frac{\bar{X}_n - \mu}{\sigma}\right)$$

**Re phrase :-**

Let $X_1, X_2, \ldots$ be i.i.d. random variables with finite mean $\mu$, finite variance $\sigma^2$. Then

$$\frac{\sqrt{n}(\bar{X} - \mu)}{\sigma} \xrightarrow{d} Z, \tag{2}$$

where $\bar{X} = \frac{X_1 + X_2 + \ldots + X_n}{n}$ and $Z \sim$ Normal $(0, 1)$.

ie. $\left| \mathbb{P}\left(\frac{\sqrt{n}(\bar{X}-\mu)}{\sigma} \leq x\right) - \mathbb{P}(Z \leq x) \right| \longrightarrow 0$

as $n \to \infty$

# Confidence Interval

$X_1, \ldots X_n$    i.i.d    $\mu$ — mean

$\sigma^2$ — variance

$$\frac{\sqrt{n}\,(\bar{X} - \mu)}{\sigma} \quad \underset{\text{Theorem}}{=}\text{Central limit} \quad \underset{}{=} \quad \begin{array}{c}\text{Normal} - Z \\ (0,1)\end{array}$$

$$\mathbb{P}\left( \frac{\sqrt{n}\,(\bar{X} - \mu)}{\sigma} \leq x \right) \quad " = " \quad \mathbb{P}(Z \leq x)$$



$$\mathbb{P}(|Z| \leq 1.96) \approx 0.95$$

$-1.96$      $1.96$

$$\bar{X} = \frac{X_1 + \ldots + X_n}{n}$$

$$\left| \frac{\sqrt{n}\,(\bar{X} - \mu)}{\sigma} \right| \quad \underset{\text{Theorem}}{=}\text{Central limit} \quad = \quad \begin{array}{c}\text{Normal} - |Z| \\ (0,1)\end{array}$$

$$\left| \frac{\sqrt{n}\,(\bar{X} - \mu)}{\sigma} \right| \leq 1.96 \quad (=) \quad -1.96 \leq \frac{\sqrt{n}(\bar{X} - \mu)}{\sigma} \leq 1.96$$

$$(=) \quad (-1.96)\sigma \leq \sqrt{n}\,(\bar{X} - \mu) \leq (1.96)\sigma$$

$$(=) \quad \bar{X} - \frac{1.96\,\sigma}{\sqrt{n}} \leq \mu \leq \bar{X} + \frac{1.96\,\sigma}{\sqrt{n}}$$

# Confidence Intervals

Assure $c \, b \, \partial - k$
$E(X) = \mu$
$\sigma$ - known

Compute :- $\bar{X} = \sum\limits_{i=1}^{n} \dfrac{x_i}{n}$

Using the Central Limit Theorem for large $n$ we have

$$P(| \frac{\sqrt{n}(\bar{X} - \mu)}{\sigma} | \leq 1.96) \approx 0.95$$

which is the same as saying

$$P(\mu \in \left( -\frac{1.96\sigma}{\sqrt{n}} + \bar{X}, \frac{1.96\sigma}{\sqrt{n}} + \bar{X} \right)) \approx 0.95$$

The interval $\left( -\frac{1.96\sigma}{\sqrt{n}} + \bar{X}, \frac{1.96\sigma}{\sqrt{n}} + \bar{X} \right)$ is called the 95% confidence interval for $\mu$.

$\hookrightarrow$ dependent on sample.
and is valid if $\sigma$ is known.

# Confidence Intervals

95% confidence interval for $\mu$ is $\left(-\frac{1.96\sigma}{\sqrt{n}} + \bar{X}, \frac{1.96\sigma}{\sqrt{n}} + \bar{X}\right)$

Meaning: for $n$ large if we did $m$ (large) repeated trials and computed the above interval for each trial then true mean would belong to approximately 95% of $m$ intervals calculated.

# Confidence Intervals

The below is code for finding the confidence interval for a data $x$.

```
> cifn = function(x, alpha=0.95){
+ z = qnorm( (1-alpha)/2, lower.tail=FALSE)
+ sdx = sqrt(1/length(x))
+ c(mean(x) - z*sdx, mean(x) + z*sdx)
+ }
```

## Three Confidence Intervals for Normal(0,1)

```
> x1 = rnorm(100,0,1);y = cifn(x1)
> y
[1] -0.35705304  0.03493976
> x2 = rnorm(100,0,1);z = cifn(x2)
> z
[1] -0.2832489   0.1087439
> x3 = rnorm(100,0,1);w = cifn(x3)
> w
[1] -0.30294682  0.08904598
```
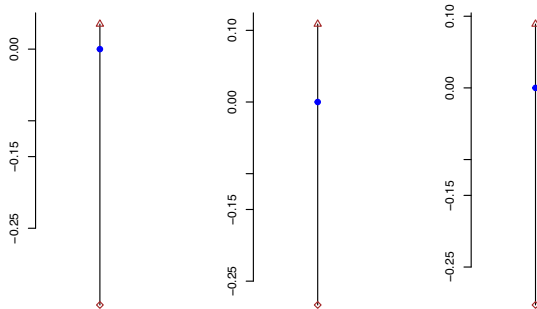
Does 0 belong to all the three confidence intervals ?

The below is a plot of the three confidence intervals computed in the previous slide.

We generate 10 trials of 100 samples from Normal(0,1) and
compute the confidence intervals using the function defined earlier.

```
> normaldata = replicate(10, rnorm(100,0,1),
+ simplify=FALSE)
> cidata = sapply(normaldata, cifn)
```

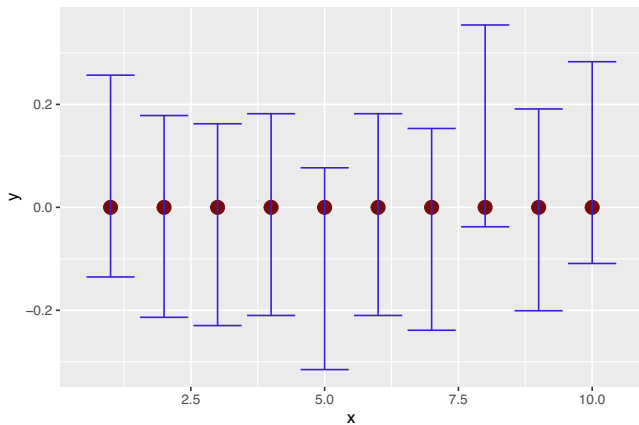It is easy to check how many of them contain 0.

```
> TRUEIN = cidata[1,]*cidata[2,]<0
> table(TRUEIN)
TRUEIN
TRUE
  10
```

# Confidence Intervals : 10 Trials

We generate 10 trials of 100 samples from Normal(0,1) and compute the confidence intervals using the function defined earlier.

```
> normaldata = replicate(40, rnorm(100,0,1),
+ simplify=FALSE)
> cidata = sapply(normaldata, cifn)
```
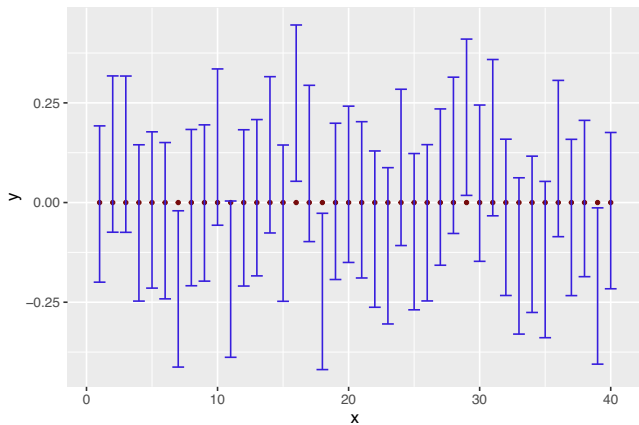
It is easy to check how many of them contain 0.

```
> TRUEIN = cidata[1,]*cidata[2,]<0
> table(TRUEIN)
TRUEIN
FALSE   TRUE
    5     35
```

We generate 100 trials of 100 samples from Normal(0,1) and compute the confidence intervals using the function defined earlier.

```
> normaldata = replicate(100, rnorm(100,0,1),
+ simplify=FALSE)
> cidata = sapply(normaldata, cifn)
```
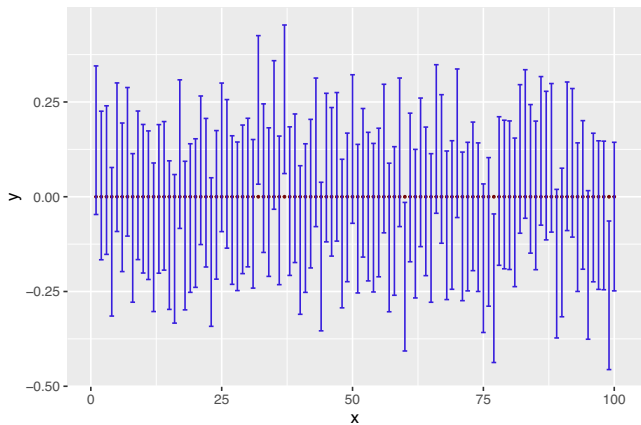
It is easy to check how many of them contain 0.

```
> TRUEIN = cidata[1,]*cidata[2,]<0
> table(TRUEIN)
TRUEIN
FALSE  TRUE
    5    95
```

# Confidence Intervals : 1000 Trials

We generate 1000 trials of 100 samples from Normal(0,1) and compute the confidence intervals using the function defined earlier.

```
> normaldata = replicate(1000, rnorm(100,0,1),
+ simplify=FALSE)
> cidata = sapply(normaldata, cifn)
```
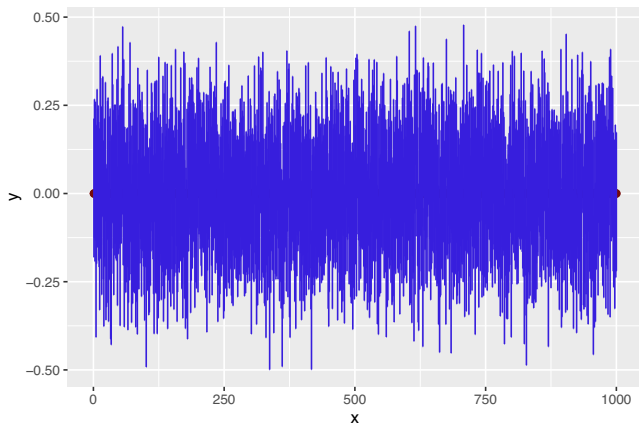
It is easy to check how many of them contain 0.

```
> TRUEIN = cidata[1,]*cidata[2,]<0
> table(TRUEIN)
TRUEIN
FALSE  TRUE
   54   946
```

95% confidence interval for $\mu$ is $\left(-\frac{1.96\sigma}{\sqrt{n}} + \bar{X}, \frac{1.96\sigma}{\sqrt{n}} + \bar{X}\right)$

Meaning: for $n$ large if we did $m$ (large) repeated trials and computed the above interval for each trial then true mean would belong to approximately 95% of $m$ intervals calculated.

Thus numerically the above meaning seems to hold for a Normal population.