

Central limit theorem & Normal distribution

$Z \sim N(0,1)$ Z is said to have normal distribution with mean 0 and variance 1 if

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

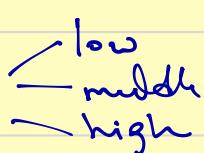
[numerical evaluation by normal tables] [improve Riemann integral]

- The distribution seems to occur naturally in many biological observations.

e.g. - Height of individuals

- # of leaves on a tree

arise as sum of independent random process

thus some contributing to 

wielding a frequency plot that is symmetrical

- This is a fact that can be established theoretically.

[De Moivre - Laplace C.L.T.]: Let Z_1, Z_2, \dots, Z_n be

i.e. $\text{Bernoulli}(p)$, $0 < p < 1$ then

$$\hat{p}_n = \frac{\sum_{i=1}^n Z_i}{n} \quad - \text{sample proportion}$$

$$\sqrt{n} \left(\hat{p}_n - p \right) \xrightarrow{d} N(0,1) \quad - \textcircled{*}$$

[i.e. Suppose $W_n = \frac{\sqrt{n}(\hat{p}_n - p)}{\sigma}$ & $Z \sim N(0,1)$

$P(W_n \leq x) \rightarrow P(Z \leq x) \text{ as } n \rightarrow \infty$]
for all $x \in \mathbb{R}$

Proof :- Remark $\textcircled{*}$:- if $S_n = \sum_{i=1}^n Z_i$
[IDEA] $\frac{S_n - np}{\sqrt{n}\sigma} \xrightarrow{d} N(0,1)$

Now, $S_n \stackrel{d}{=} \text{Binomial}(n, p)$

$$P(a \leq \frac{S_n - np}{\sqrt{n}\sigma} < b) = \int_a^b \frac{e^{-y^2/2}}{\sqrt{\pi}} dy$$

$$\sigma = \sqrt{np(1-p)}$$

$$q = 1-p \quad | \quad \sum_{k=[a\sqrt{np(1-p)} + np]}^{[b\sqrt{np(1-p)} - 1 + np]} P(S_n = k) - \int_a^b \frac{e^{-y^2/2}}{\sqrt{\pi}} dy$$

$$\begin{aligned}
 &\leq \left| \sum_{k=[a\sqrt{n}+\epsilon n p]}^{[b\sqrt{n}+\epsilon n p]} \left[P(S_n=k) - \frac{1}{\sqrt{n} p \sigma} \phi\left(\frac{k-np}{\sqrt{n} p \sigma}\right) \right] \right| \\
 &+ \left| \sum_{k=[a\sqrt{n}+\epsilon n p]}^{[b\sqrt{n}+\epsilon n p]} \frac{1}{\sqrt{n} p \sigma} \phi\left(\frac{k-np}{\sqrt{n} p \sigma}\right) - \int_a^b \frac{e^{-\frac{s^2}{2}}}{\sqrt{2\pi}} ds \right| \\
 &= \text{I} + \text{II}
 \end{aligned}$$

II - Riemann Sum Convergence

$$\begin{aligned}
 \text{I} - P(S_n=k) &= \frac{n!}{k!(n-k)!} b^k a^{n-k} \\
 &\bullet n! \sim n^n e^{-n} \sqrt{n} \cdot \sqrt{2\pi} \quad [\text{Using Stirling's formula}]
 \end{aligned}$$

$$\text{I} + \text{II} \rightarrow 0 \text{ as } n \rightarrow \infty \quad \square$$

Central Limit Theorem: X_1, \dots, X_n are i.i.d random variables with $E X_i = \mu$ & $\text{Var}(X_i) = \sigma^2$

Then

$$\sqrt{n} \left(\overrightarrow{\bar{X}} - \mu \right) \xrightarrow{d} N(0,1).$$

Proof :- uses several results ... \square
beyond the scope of this class

Central Limit Theorem

Suppose we want to verify the below result via simulations:

Let X_1, X_2, \dots be i.i.d. random variables with finite mean μ , finite variance σ^2 . Then

$$\frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \xrightarrow{d} Z, \quad (1)$$

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

Central Limit Theorem

Suppose we want to verify the below result via simulations:

Let X_1, X_2, \dots be i.i.d. random variables with finite mean μ , finite variance σ^2 . Then

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

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

Central Limit Theorem

We could rephrase the result as:

Let X_1, X_2, \dots be i.i.d. random variables with finite mean μ , finite variance σ^2 . Then

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

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

Central Limit Theorem

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

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

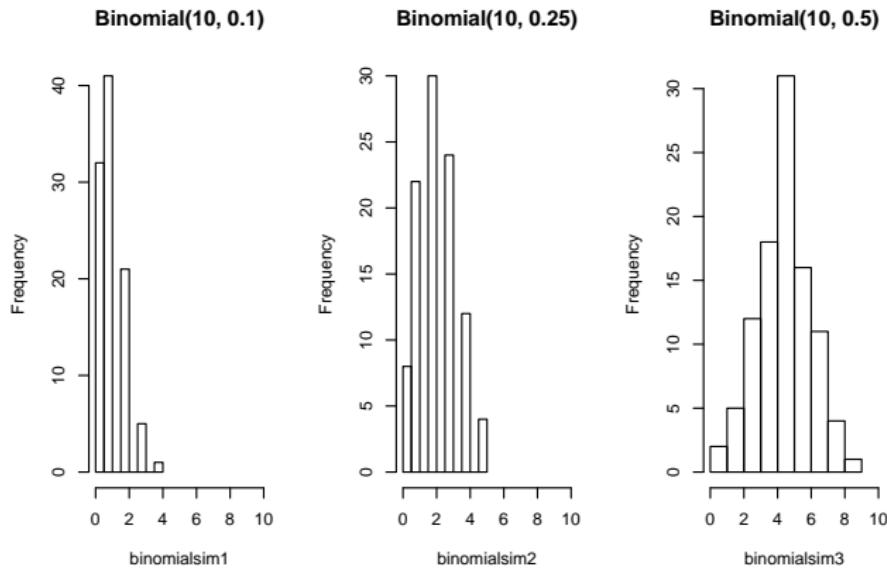
is close to a Normal distribution.

Central Limit Theorem

We may simulate Binomial samples either directly by `rbinom` command or using the `replicate` and `rbinom` command.

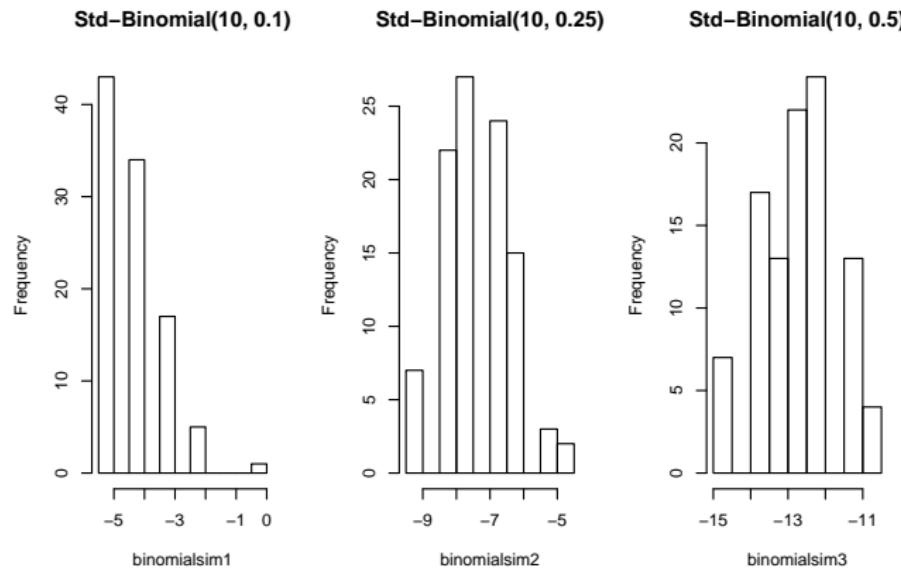
```
> binomialsim1 = rbinom(100,10,0.1)
> # generates 100 Binomial (10,0.1) samples
>
> binomialsim2 = replicate(100, rbinom(1,10,0.25))
> # generates 100 Binomial (10,0.25) samples
>
> binomialsim3 = replicate(100, rbinom(1,10,0.5))
> # generates 100 Binomial (10,0.5) samples
>
```

Histogram of all three simulations



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

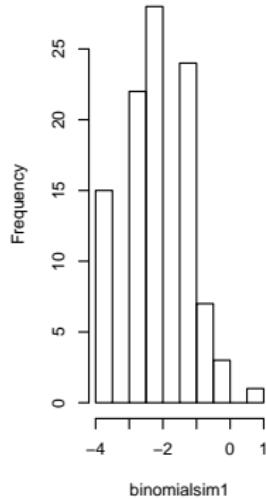
Standardised Histograms: Binomial $n=10$ and $p=0.1, 0.25, 0.5$



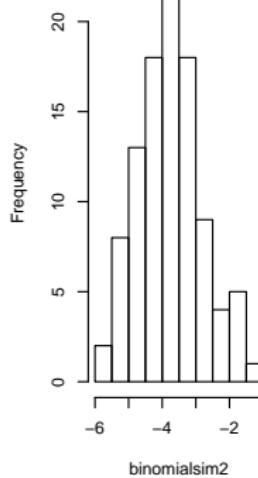
Perhaps $n = 10$ is not large enough to see the Central Limit Theorem occurring.

Standardised Histograms: Binomial $n=20$ and $p=0.1, 0.25, 0.5$

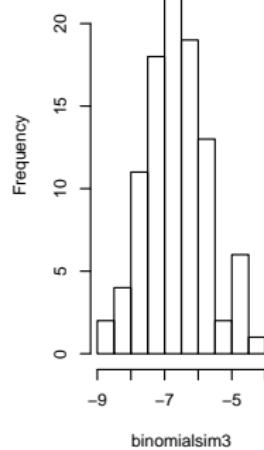
Std-Binomial(20, 0.1)



Std-Binomial(20, 0.25)

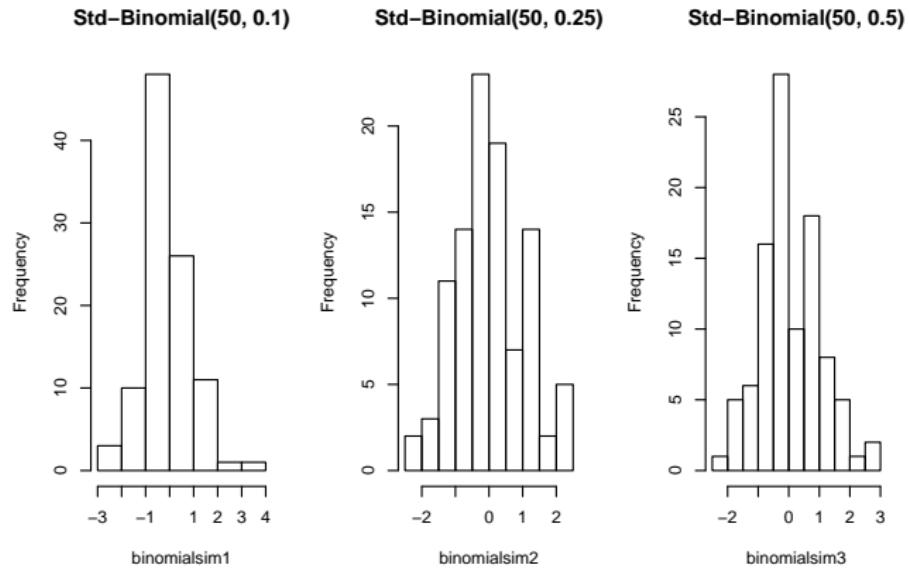


Std-Binomial(20, 0.5)



$n = 20$ is better.

Standardised Histograms: Binomial $n=50$ and $p=0.1, 0.25, 0.5$



$n = 50$ 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.