

# 14) Central Limit Theorem

Theorem: 4.1 Let  $X \sim N(\mu_X, \sigma_X^2)$   $Y \sim N(\mu_Y, \sigma_Y^2)$  be two independent random variables with  $\mu_X, \sigma_X \in \mathbb{R}$  and  $\mu_Y, \sigma_Y \in \mathbb{R}$ . Then

$$X + Y \sim N(\mu_X + \mu_Y, \sigma_X^2 + \sigma_Y^2)$$

proof

We know that for normal distribution

$$E[X] = \mu_X \quad \text{Var}(X) = \sigma_X^2$$

$$E[Y] = \mu_Y \quad \text{Var}(Y) = \sigma_Y^2$$

$$E[X + Y] = E[X] + E[Y]$$

$$= \mu_X + \mu_Y$$

linearity of  
expectations  
Thm 10.2

$$\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + \text{Cov}(X, Y)$$

$$X \perp Y \Rightarrow \text{Cov}(X, Y) = 0. \text{ So}$$

$$\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y)$$

$$= \sigma_X^2 + \sigma_Y^2.$$

$$\text{Hence } X + Y \sim N(\mu_X + \mu_Y, \sigma_X^2 + \sigma_Y^2)$$

↳ why normal, see year 2 [?]

Now consider an iid sample  $x_1, \dots, x_n$  from a random variable  $X$

Sample mean

$$\bar{X}_n = \frac{(x_1 + x_2 + \dots + x_n)}{n}$$

From chapter 13:

$$E[\bar{X}_n] = E[X]$$

$$\text{Var}(\bar{X}_n) = \frac{\text{Var}(X)}{n}$$

As seen in R labs, for large sample size, the sample mean is normally distributed (approximately)

↳ even if  $X$  is not normally distributed.

Rule of thumb  
14.2 For sufficiently large sample size  $n$ , mean is approximately distributed,

$$\bar{X}_n = (X_1 + X_2 + \dots + X_n)/n \dot{\sim} N\left(E[X], \frac{\text{Var}(X)}{n}\right)$$

Equivalently multiplying sample mean  $\bar{X}_n$  by  $n$  gives:

$$n\bar{X}_n = X_1 + X_2 + X_3 + \dots + X_n$$

$$E[X_1 + X_2 + \dots + X_n] = E[n\bar{X}_n]$$

$$= nE[\bar{X}_n]$$

$$= nE[X] \quad \text{Thm 7.16}$$

$$\text{Var}(X_1 + X_2 + \dots + X_n) = \text{Var}(n\bar{X}_n)$$

$$= n^2 \text{Var}(\bar{X}_n) \quad \text{Thm 7.25}$$

$$= n^2 \frac{\text{Var}(X)}{n}$$

$$= n \text{Var}(X)$$

So

by Thm  
14.1

$$X_1 + X_2 + \dots + X_n \dot{\sim} N(nE[X], n\text{Var}(X))$$

The amazing fact is that no matter what the distribution of the random variable  $X$  is, for sufficiently large sample size  $n$ , sample mean is approximately normally distributed.

14.2 is not a theorem as the term sufficiently large is vague.

To get a precise theorem, we need to take limit as  $n \rightarrow \infty$ .

$X_n$  does not have a nice limit as  $n \rightarrow \infty$  as its variance goes towards 0.

Therefore we consider the standardised random variable

$$Z_n = \frac{\bar{X}_n - E[\bar{X}_n]}{\sqrt{\text{Var}(\bar{X}_n)}} = \frac{\bar{X}_n - E[X_i]}{\sqrt{\frac{\text{Var}(X_i)}{n}}} \approx N(0,1)$$

which for all  $n$  has variance one

proof

$$E[Z_n] = E \left[ \frac{\bar{X}_n - E[\bar{X}_n]}{\sqrt{\text{Var}(\bar{X}_n)}} \right]$$

linearity  
of  
expectations

numbers  
in  $\mathbb{R}$

$$= E \left[ \frac{\bar{X}_n}{\sqrt{\text{Var}(\bar{X}_n)}} \right] - E \left[ \frac{E[\bar{X}_n]}{\sqrt{\text{Var}(\bar{X}_n)}} \right]$$

$$= \frac{E[\bar{X}_n]}{\sqrt{\text{Var}(\bar{X}_n)}} - \frac{E[\bar{X}_n]}{\sqrt{\text{Var}(\bar{X}_n)}}$$

$$= 0 \Rightarrow E[Z_n] = 0$$

$$\text{Var}(Z_n) = \text{Var} \left( \frac{\bar{X}_n - E[\bar{X}_n]}{\sqrt{\text{Var}(\bar{X}_n)}} \right)$$

$$= \text{Var} \left( \frac{\bar{X}_n}{\sqrt{\text{Var}(\bar{X}_n)}} - \frac{E[\bar{X}_n]}{\sqrt{\text{Var}(\bar{X}_n)}} \right)$$

by Thm 7.25

$$= \frac{\text{Var}(\bar{X}_n)}{\text{Var}(\bar{X}_n)} = 1 \Rightarrow \text{Var}(Z_n) = 1$$

So

$$Z_n \sim N(0, 1)$$

Theorem: (Central Limit Theorem):  
14.3

For any  $n \in \mathbb{N}$ , let  $X_1, X_2, \dots$  be an iid sample from a distribution with finite expectation  $\mu$  and finite variance  $\sigma^2$ .

Let

$$Z_n = \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}}$$

Then at any point  $x \in \mathbb{R}$

$$\lim_{n \rightarrow \infty} F_{Z_n}(x) = \phi(x)$$

Convergence in distribution, where  $\phi$  is the distribution function of standard normal distribution.

Example: (Using smarties example, chapter 10):

$Y$ : no of yellow smarties

$$Y = \sum_{i=1}^n y_i \quad y_i = \begin{cases} 1 & \text{if smartie is yellow} \\ 0 & \text{otherwise} \end{cases}$$

$$[n=40, p_Y = \frac{1}{8}]$$

$$\bar{y}_n = \frac{1}{n} \sum_{i=1}^n y_i$$

$P(|\bar{Y}_n - p_Y| > 0.1) \rightarrow$  estimate that probability that difference of estimate  $\bar{Y}_n$  and true value  $p_Y = 1/8$  is greater than 0.1

$\hookrightarrow$  "probability that we are 0.1 away from our estimate"

Since  $y_i$  is an indicator random variable, it is a Bernoulli distributed so

$$E[y_i] = \mu = p_Y \quad \left( \begin{array}{l} \text{for } X \sim \text{Ber}(p) \\ E[X] = p \end{array} \right)$$

For bernoulli

$$\text{Var}(Y_i) = p(1-p) > 0$$

$$P(|\bar{Y}_n - p_Y| > 0.1) = P\left(\frac{|\bar{Y}_n - p_Y|}{\sqrt{\frac{\text{Var}(Y_i)}{n}}} > \frac{0.1}{\sqrt{\frac{\text{Var}(Y_i)}{n}}}\right)$$

$$= P\left(\frac{|\bar{Y}_n - p_Y|}{\sqrt{\frac{(1-p)p}{n}}} > \frac{0.1}{\sqrt{\frac{(1-p)p}{n}}}\right)$$

$$\text{let } Z_n = \frac{\bar{Y}_n - E[Y_i]}{\sqrt{\frac{\text{Var}(Y_i)}{n}}} = \frac{\bar{Y}_n - p_Y}{\sqrt{\frac{(1-p)p}{n}}} \sim N(0,1)$$

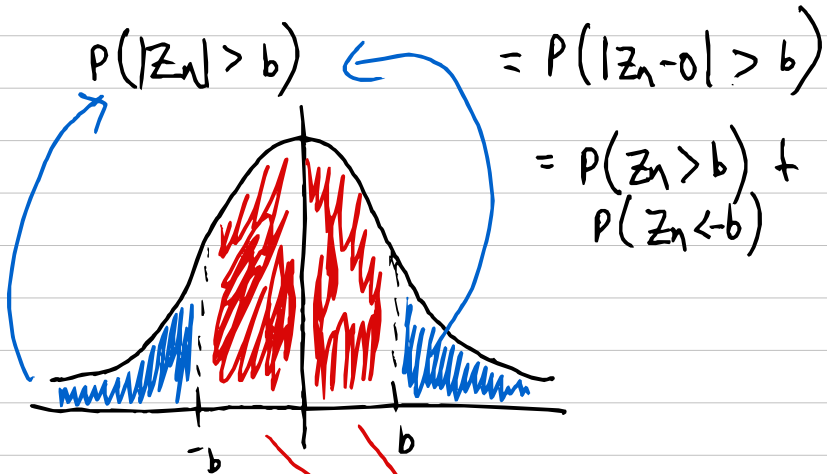
$$\text{let } \frac{0.1}{\sqrt{\frac{(1-p)p}{n}}} = b$$

$$= P(|Z_n| > b)$$



Need to calculate

$$\left( \begin{array}{l} P(Z_n > b) \text{ or } (+) \\ P(-Z_n > b) \\ \hline P(Z_n > b) + P(Z_n < -b) \end{array} \right)$$



or

$$\begin{aligned} P(|Z_n| > b) &= 1 - P(|Z_n| \leq b) \\ &= 1 - P(-b \leq Z_n \leq b) \\ &= 1 - (P(Z_n \leq b) - P(Z_n \leq -b)) \\ &= 1 - P(Z_n \leq b) + P(Z_n \leq -b) \\ &= P(Z_n > b) + P(Z_n < -b) \end{aligned}$$

Due to symmetry of normal distribution,

$$P(Z_n > b) = P(Z_n < -b) \quad \left( \begin{array}{l} \text{check week 5} \\ \text{exercises for proof} \end{array} \right)$$

So

$$P(|Z_n| > b) = P(Z_n > b) + P(Z_n < -b)$$

$$= 2P(Z_n < -b)$$

$$= 2\Phi(-b)$$

↳ distribution fn of  
standard  
normal distribution

$$\approx 0.06$$

Calculating actual value, not approximation  
using central limit thm:

$$Y \sim \text{Bin}(n, p_Y) \Rightarrow Y \sim \text{Bin}(40, 1/8)$$

$$P(|\bar{Y}_n - p_Y| > 0.1) = P(n|\bar{Y}_n - p_Y| > 0.1n)$$

$$\text{(since } n > 0) \quad = P(|n\bar{Y}_n - np_Y| > 0.1n)$$

$$\left[ \begin{array}{l} n\bar{Y}_n = Y \text{ as} \\ \bar{Y}_n = \frac{1}{n} \sum_{i=1}^n x_i, \\ Y = \sum_{i=1}^n x_i \end{array} \right] = P(|Y - 5| > 4)$$

$$P(|Y-S| > 4)$$

$$= 1 - P(|Y-S| \leq 4)$$

$$= 1 - [P(-4 \leq Y-S \leq 4)]$$

$$= 1 - (P(Y-S \leq 4) - P(Y-S \leq -4))$$

$$= 1 - P(Y \leq 9) + P(Y \leq 1)$$

$$= 0.05487806$$

$$\approx 0.06$$

pretty accurate to approximation by  
central limit thm.

→ used `pbinom(q, n, p)`  
in R