

## Estimating $p = P(X \in A)$ for a given region $A$ :

Simulate a large number ( $N$ ) of independent draws from the distribution of  $X$ , say,  $X_1, X_2, \dots, X_N$

Define  $Y_1, \dots, Y_N$  as:

$$Y_i = I(X_i \in A)$$

$\Rightarrow Y_1, Y_2, \dots, Y_N$  are draws from Bernoulli( $p$ )

MC estimator of  $p$ :

$$\hat{p} \approx \frac{1}{N} \sum_{i=1}^N Y_i$$

Properties of  $\hat{p}$ :

/ see next page

Note:

Define:

$$Y = I(X \in A)$$

$$Y = \begin{cases} 1, & \text{if } X \in A \\ 0, & \text{o/w.} \end{cases}$$

Then:  $Y \sim \text{Bernoulli}(p)$

$$P(X \in A) = p$$

$$E[Y] = p$$

$$\text{Var}[Y] = p(1-p)$$

Properties of  $\bar{X}$  [based on a large sample of size  $N$ ]:

• LLN: If  $N$  is large,  $\bar{X} \approx \mu = E(X)$

• CLT: If  $N$  is large,  $\bar{X} \sim N[\mu = E(X), \frac{\sigma^2}{N}]$

$$\downarrow$$
$$\equiv Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{N}} = \frac{\bar{X} - \mu}{\sigma/\sqrt{N}} \sim N(0,1)$$

Properties of  $\hat{p}$  [based on a large sample of size  $N$ ]

• LLN: If  $N$  is large,  $\hat{p} \approx p$

• CLT: If  $N$  is large,  $\hat{p} \sim N[p, \frac{p(1-p)}{N}]$

$$\equiv Z = \frac{\hat{p} - p}{\sqrt{p(1-p)/N}} \sim N(0,1)$$

# Accuracy of a Monte Carlo study: (for estimation of $p$ ):

Error in estimation:  $\hat{p} - p$

Specify a small margin of error  $\epsilon$  and a small probability  $\alpha$ .

Want ( $N$ ) such that

$$P(|\hat{p} - p| > \epsilon) \leq \alpha \quad (1)$$

or equivalently

$$P(|\hat{p} - p| \leq \epsilon) \geq 1 - \alpha$$

need to calculate using CLT.

Need to solve this for  $N$ .

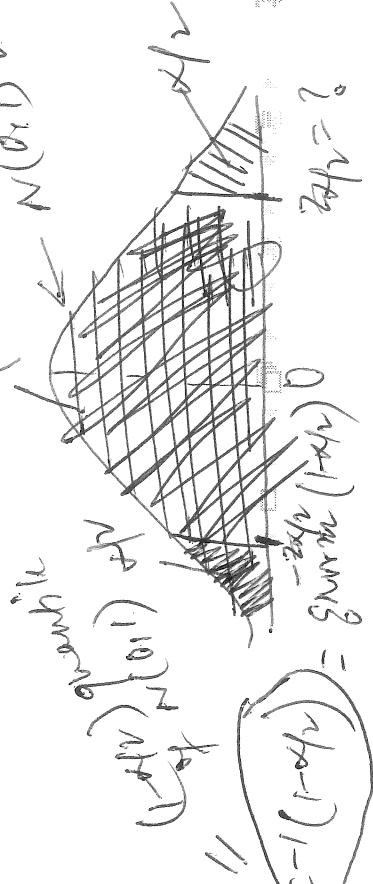
- Error exceeds  $\epsilon$  with probability  $\alpha$  or less
- Error is  $\epsilon$  or less with probability more than  $1 - \alpha$ .

To derive a formula for  $N$ , suppose  $Z \sim N(0, 1)$  and  $z_{\alpha/2}$  is such that  $P(Z > z_{\alpha/2}) = \alpha/2$ .

$$P[Z \leq z_{\alpha/2}] = 1 - \alpha/2$$

$$\Rightarrow F(z_{\alpha/2}) = 1 - \alpha/2$$

$$\text{cdf of } N(0,1) \Rightarrow z_{\alpha/2} = F^{-1}(1 - \alpha/2)$$



By def:  $P[|Z| > z_{\alpha/2}] = \alpha \equiv P[|Z| \leq z_{\alpha/2}] = 1 - \alpha$

From the symmetry,

$$P(|Z| > z_{\alpha/2}) = \alpha \quad (2)$$

Now, let's derive an expression for  $P(|\hat{p} - p| > \epsilon)$ :

$$= P\left[\frac{|\hat{p} - p|}{\sqrt{\frac{p(1-p)}{N}}} > \frac{\epsilon}{\sqrt{\frac{p(1-p)}{N}}}\right] \approx P\left[|Z| > \frac{\epsilon}{\sqrt{p(1-p)/N}}\right] \leq \alpha = P[|Z| > z_{\alpha/2}]$$

using CLT.

$$\hat{p} - p \sim N(0, 1)$$

Comparing (2) and (3), and noticing that  $P(|Z| > x)$  is decreasing in  $x$ , we can conclude that (1) approximately holds if

$$\frac{\epsilon}{\sqrt{\frac{p(1-p)}{N}}} \geq z_{\alpha/2} \Rightarrow \frac{\epsilon^2}{\frac{p(1-p)}{N}} \geq z_{\alpha/2}^2 \Rightarrow \frac{N\epsilon^2}{p(1-p)} \geq z_{\alpha/2}^2$$

$$\Rightarrow N \geq \frac{z_{\alpha/2}^2 p(1-p)}{\epsilon^2}$$

Take  $N = \frac{z_{\alpha/2}^2 p(1-p)}{\epsilon^2}$

A practical problem:

The answer depends on  $p$ , which is unknown.

Alternative 1:

Replace  $p$  by a 'good' guess, say  $p^*$ .

Alternative 2:

Replace  $p(1-p)$  by its maximum possible value, to get:

$$N \approx \frac{2^2 / 4}{e^2}$$

sample size formula to use in practice.



**Note:** This formula is valid only if  $N$  is large.

Ex: Suppose the desired accuracy is  $(\epsilon, \alpha) = (0.03, 0.05)$ .  $N = ?$

$\alpha/2 = 0.025$   
 $\Rightarrow Z_{0.025} = \text{norm}(1-0.025)$   
 $\approx 1.96 \rightarrow R$

$P[| \text{Error} | \leq 0.03] \geq 0.95$

choice  
 common practice  
 in polling  
 for companies

$$N \approx \frac{(1.96)^2 \cdot \frac{1}{4}}{(0.03)^2} = 1067.11$$

rounding up, give 1068

$N \approx 1068$