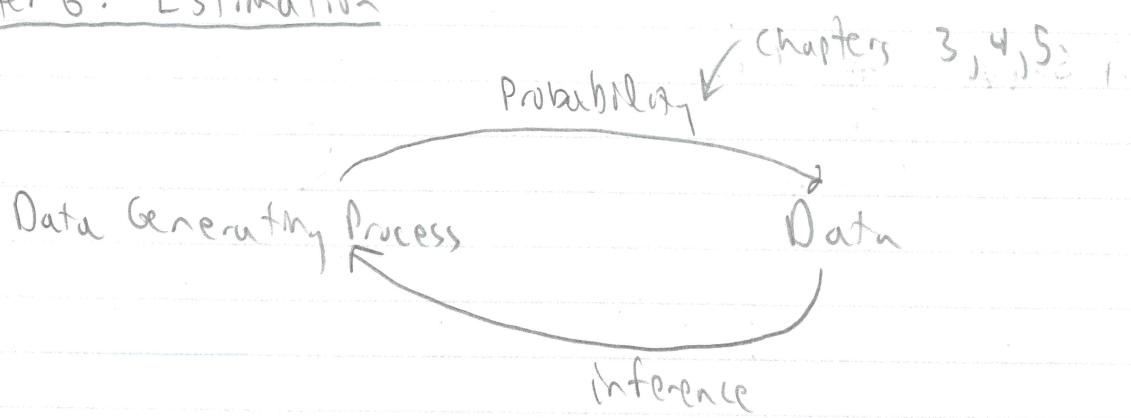


## Chapter 6: Estimation



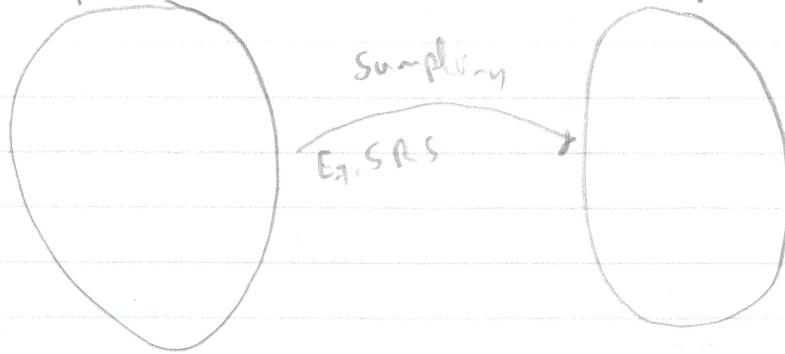
- Probability:  $X \sim N(80, 12)$   
    ↑ assumed know  
    What is  $\Pr(X > 90)$
- Inference: Observe  $X_i = 81, 78, 77, 89, \dots$   
    ↑ Assume  $X_i \sim N(\mu, \sigma^2)$  what are  $\mu$  and  $\sigma^2$ ?
- Estimation: Guess parameter values from data
  - Point estimation: A single number is your best guess  
    ↑ E.g. estimate  $\mu$  with  $\bar{X}$
  - Interval estimation: Get a range of likely values of a parameter
- Confidence Intervals
- Hypothesis Testing: How sure are we a parameter is different from some value?  
    ↑ e.g. non-zero?

group we are interested in

group we have data about

Reference / target / study population

Sample



parameter: summary of pop.

$$\mu, \sigma^2, p$$

statistic: summary of sample

$$\bar{x}, s^2, \hat{p}$$

- Simplest way to get a sample is by a simple random sample

Each unit has an equal chance of being in sample

- Random selection (via SRS) is distinct from

Random Assignment: Randomly assign units to different groups (e.g. treatment vs control)

• Random Selection: results generalizable to target population (b/c sample similar to pop in terms of demographic variables, etc.)

• Random Assignment: Allows for claims of causality because all possible confounders equal in the groups.

- Randomized Clinical Trial (RCT): Random assignment of treatment to compare them
- No causal claims without random assignment

E.g. tobramycin and gentamicin are antibiotics, tobramycin is more aggressive and has more side effects. Early studies were not randomized and showed tobramycin performed worse. Why?

↑ Doctors gave sicker patients tobramycin

↑ Randomization guarantees equal # sicker and less sick in each group (on average).

- Double blind - neither doctor nor patient know treatment

↑ guards against placebo effect

- Single blind - doctor knows

- unblinded - both know

Sampling in R

- Estimate mean:

Suppose  $E[X_i] = \mu$

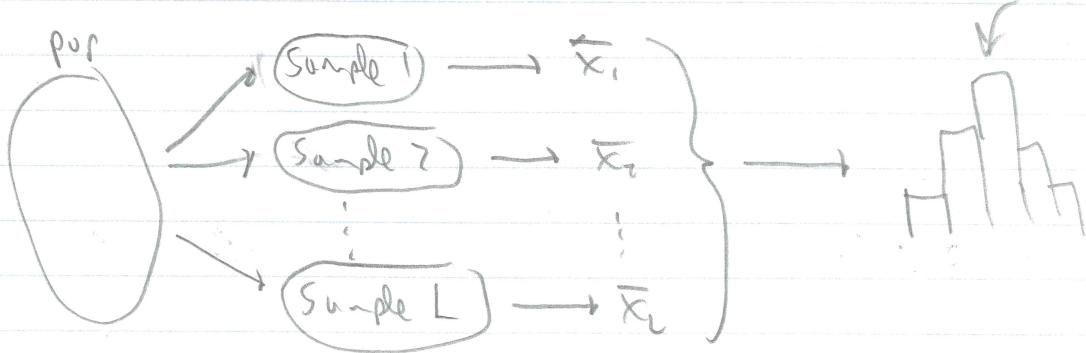
↑ Not necessarily normal

Estimate  $\mu$  with  $\hat{\mu} = \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$

↑ estimate of  $\mu$  [sample mean]

- Sampling distribution: Distribution of statistic across many (hypothetical) samples.

sampling distribution



↑ Used to describe properties of statistics

- Properties of  $\bar{X}$

Let  $X_1, X_2, \dots, X_n$  be a random sample with  
 $E(X_i) = \mu, \text{Var}(X_i) = \sigma^2$

The -

- $E(\bar{X}) = \mu$  (unbiased)

- $\text{Var}(\bar{X}) = \sigma^2/n$

more precise for larger  $n$  (consistent)

3)  $\bar{X} \approx N(\mu, \sigma^2/n)$  for large  $n$ , even if  $X_i$  are not also normal  
 (central limit theorem)

- Standard error: Standard deviation of a statistic

$$SE(\bar{X}) = \sigma / \sqrt{n}$$

↑  
not individual observations

- Estimated Standard Error (Standard Error): estimated standard deviation of a statistic

Book just means estimated SE when they say SE, typically

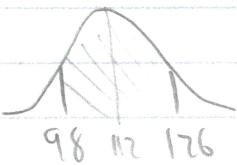
$$SE(\bar{X}) = s / \sqrt{n}$$

- typically don't know  $\sigma^2$ , so estimate it with  $s^2$

<u>Statistic</u>	$X_i$	$\bar{X}$
Mean	$\mu$	$\mu$
Variance	$\sigma^2$	$\sigma^2/n$
Estimated Variance	$s^2$	$s^2/n$
Distribution	Unknown	$N(\mu, \sigma^2/n)$ (for large $n$ )

- Exercise: Mean birthweight is 112 oz with standard deviation 20.6, what is the probability the mean of 10 birthweights will be between 98 and 126?

$$\bar{X} \sim N(112, (20.6/\sqrt{10})) = N(112, 6.514^2)$$



$$\begin{aligned}
 &= \text{pnorm}(126, 112, 6.514) - \text{pnorm}(98, 112, 6.514) \\
 &= 0.9684
 \end{aligned}$$

## • Interval Estimation

- $\bar{X}$  is not exactly equal to  $p$
- Want range of likely values of  $p$

• Know  $\bar{X} \sim N(p, \sigma^2/n)$  for large  $n$

$$\Rightarrow \Pr(-1.96 \leq \frac{\bar{X}-p}{\sigma/\sqrt{n}} \leq 1.96) \approx 0.95$$

$$\Rightarrow \Pr\left(-1.96 \frac{\sigma}{\sqrt{n}} \leq \bar{X}-p \leq 1.96 \frac{\sigma}{\sqrt{n}}\right) \approx 0.95$$

$$\Rightarrow \Pr\left(\bar{X} - 1.96 \frac{\sigma}{\sqrt{n}} \leq p \leq \bar{X} + 1.96 \frac{\sigma}{\sqrt{n}}\right) \approx 0.95$$

$\Rightarrow \bar{X} \pm 1.96 \frac{\sigma}{\sqrt{n}}$  captures  $p$  w.p. 0.95

↳ 95% CI

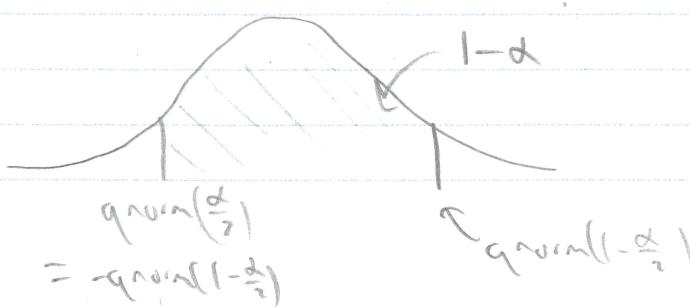
• Estimate  $\pm$  multiplier  $\times$  standard error

↑ common CI format

• More generally,

$$\boxed{\bar{X} \pm q_{\text{norm}}(1-\frac{\alpha}{2}) \frac{\sigma}{\sqrt{n}}}$$

$(1-\alpha) 100\%$  CI for  $p$



## R Interpretation of CI

- The above only works when  $\sigma^2$  is known.

↑  $\sigma^2$  is never known.

- $\frac{\bar{X} - \mu}{s/\sqrt{n}} \sim t_{n-1}$  (not  $N(0, 1)$ !)

↑ t-distribution with  $n-1$  degrees of freedom

- Only an exact result when  $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$

↑ But t-distribution tends to work better in small samples even when  $X_i$  are not normal

- For large  $n$ ,  $t_{n-1} \approx N(0, 1)$ , so CLT is ok.

- Bell shaped, centered at 0

- df ↓  $\Rightarrow$  extreme values more likely
- df ↑  $\Rightarrow$  extreme values less likely

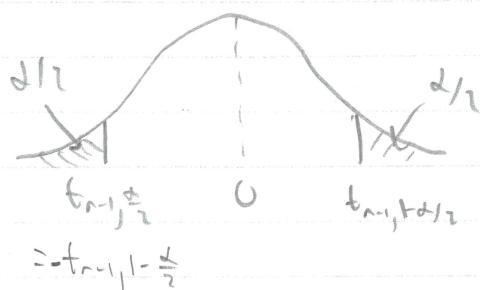
- Use t b/c of added variability from using  $s^2$  instead of  $\sigma^2$

## R code for t-distribution

- Rosner uses notation  $t_{df, p}$  for the  $p$  quantile of a  $t_{df}$  distribution.

$$t_{df, p} = qt(p, df)$$

$$\Pr\left(-t_{n-1, 1-\frac{\alpha}{2}} \leq \frac{\bar{X} - \mu}{s/\sqrt{n}} \leq t_{n-1, 1-\frac{\alpha}{2}}\right) = 1 - \alpha$$



$$\Rightarrow \Pr\left(-t_{n-1, 1-\frac{\alpha}{2}} \frac{s}{\sqrt{n}} \leq \bar{X} - \mu \leq t_{n-1, 1-\frac{\alpha}{2}} \frac{s}{\sqrt{n}}\right) = 1 - \alpha$$

$$\Rightarrow \Pr\left(\bar{X} - t_{n-1, 1-\frac{\alpha}{2}} \frac{s}{\sqrt{n}} \leq \mu \leq \bar{X} + t_{n-1, 1-\frac{\alpha}{2}} \frac{s}{\sqrt{n}}\right) = 1 - \alpha$$

$$\Rightarrow \boxed{\bar{X} \pm t_{n-1, 1-\frac{\alpha}{2}} \frac{s}{\sqrt{n}}}$$

$\uparrow (1 - \alpha) 100\% \text{ CI for } \mu$

- Exercise:  $n=10$ ,  $\bar{x} = 116.90$ ,  $s = 21.70$

calculate 90%, 95%, 99% CI's

Solution:  $qt(0.95, 9) = 1.833$

$$qt(0.975, 9) = 2.262$$

$$qt(0.995, 9) = 3.25$$

$$116.90 \pm 1.833 \cdot 21.70 / \sqrt{10}$$

$$116.90 \pm 2.262 \cdot 21.70 / \sqrt{10}$$

$$116.90 \pm 3.25 \cdot 21.70 / \sqrt{10}$$

• Note: CI level  $\uparrow$  ( $\text{so } \alpha \downarrow$ )  $\Rightarrow$  larger intervals

$n \uparrow \Rightarrow$  smaller intervals

$s^2 \uparrow \Rightarrow$  larger intervals

### Bone Density Case Study

#### Estimate Variance:

Estimate  $\sigma^2$  with  $s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$

• Why  $n-1$ ?

$$\bar{x} = \underset{a}{\operatorname{argmin}} \sum_{i=1}^n (x_i - a)^2$$

ideal

$$\Rightarrow \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \leq \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$

too small

$\Sigma p_n$  dividing by  $n-1$  makes it bigger.

$$\frac{s^2}{\sigma^2/(n-1)} \sim \chi^2_{n-1}$$

Chi-squared distribution with  $n-1$  degrees of freedom

- Properties of Chi-squared

- Support  $\geq 0$

- $E[\chi^2_{df}] = df$

- $df \downarrow \Rightarrow$  thicker tails (extreme events happen more often)

- Can use this to get confidence intervals for  $\sigma^2$

- Let  $\chi^2_{df,p} = p^{\text{th}}$  quantile of a  $\chi^2_{df}$  distribution

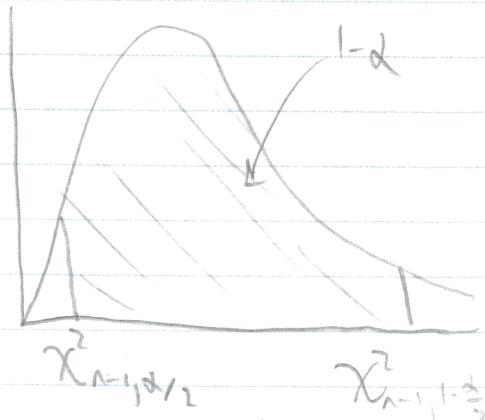
$$\Pr(\chi^2_{n-1,1-\alpha/2} \leq \frac{s^2}{\sigma^2/(n-1)} \leq \chi^2_{n-1,\alpha/2}) = 1-\alpha$$

$$\Rightarrow \Pr\left[\frac{(n-1)s^2}{\chi^2_{n-1,1-\alpha/2}} \leq \sigma^2 \leq \frac{(n-1)s^2}{\chi^2_{n-1,\alpha/2}}\right] = 1-\alpha$$

$$\frac{(n-1)s^2}{\chi^2_{n-1,\alpha/2}} \leq \sigma^2 \leq \frac{(n-1)s^2}{\chi^2_{n-1,1-\alpha/2}}$$

100(1- $\alpha$ )% CI

$\alpha$



- Notes:
  - Less often constructed than CI for  $p$
  - Very sensitive to violations of normality
  - 2 CLT does not save you

- $X_1, X_2, \dots, X_n$  be independent Bernoulli trials st.

$$X_i = \begin{cases} 1 & \text{w.p. } p \\ 0 & \text{w.p. } 1-p \end{cases}$$

Then  $X = \sum_{i=1}^n X_i = \# \text{ of } 1's \text{ out of } n \text{ trials}$

$$X \sim \text{Bin}(n, p)$$

- Example: Estimate prevalence of malignant melanoma  
Sample  $n = 5000$  individuals,  $X_i = \begin{cases} 1 & \text{if has melanoma} \\ 0 & \text{if no} \end{cases}$

- Goal: Observe  $X$ , estimate  $p$

$$\hat{p} = \frac{1}{n} \sum_{i=1}^n X_i = \frac{1}{n} X = \text{proportion of } 1's$$

$$\bullet E(\hat{p}) = E\left[\frac{1}{n} X\right] = \frac{1}{n} E(X) \stackrel{\text{Bernoulli r.v.}}{=} \frac{1}{n} np = p \Rightarrow \text{unbiased}$$

$$\bullet \text{Var}(\hat{p}) = \text{Var}\left(\frac{1}{n} X\right) = \frac{1}{n^2} \text{Var}(X) = \frac{1}{n^2} n p(1-p) = \frac{p(1-p)}{n}$$

↑ More precise for larger  $n \Rightarrow$  consistent

- Estimated standard error:  $\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$
  - Goal: Interval estimate of  $p$
  - For large  $n$ ,  $\hat{p} \sim N(p, \frac{p(1-p)}{n})$   
 $n\hat{p}(1-\hat{p}) \geq 5$
  - let  $Z_p = q_{\text{norm}}(p) = p^{\pm}$  quantile of  $N(0, 1)$
  - $\frac{\hat{p} - p}{\sqrt{\hat{p}(1-\hat{p})/n}} \approx N(0, 1)$
  - $\Rightarrow \Pr(-Z_{1-\alpha/2} \leq \frac{\hat{p} - p}{\sqrt{\hat{p}(1-\hat{p})/n}} \leq Z_{1-\alpha/2}) \approx 0.95$
- Solve for  $p$
- $$\hat{p} \pm Z_{1-\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$
- is a  $(1-\alpha)100\%$  CI for  $p$
- Estimate      Multiplier      Standard error

- Exercise: 1000 women, 400 have breast cancer. What is a 95% CI for breast cancer incidence?

Solution:  $\hat{p} = 0.4$ ,  $n\hat{p}(1-\hat{p}) = 384 > 5$

$$Z_{0.975} = q_{\text{norm}}(0.975) \approx 1.96$$

$$0.4 \pm 1.96 \cdot \sqrt{0.4 \cdot (1-0.4)/1000} = (0.03616, 0.04384)$$

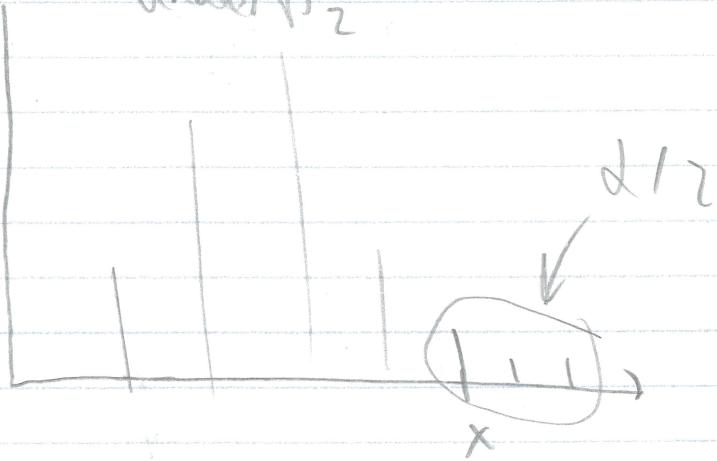
• What if  $n$  is small? Use an exact method

Find  $p_1$  and  $p_2$  s.t.

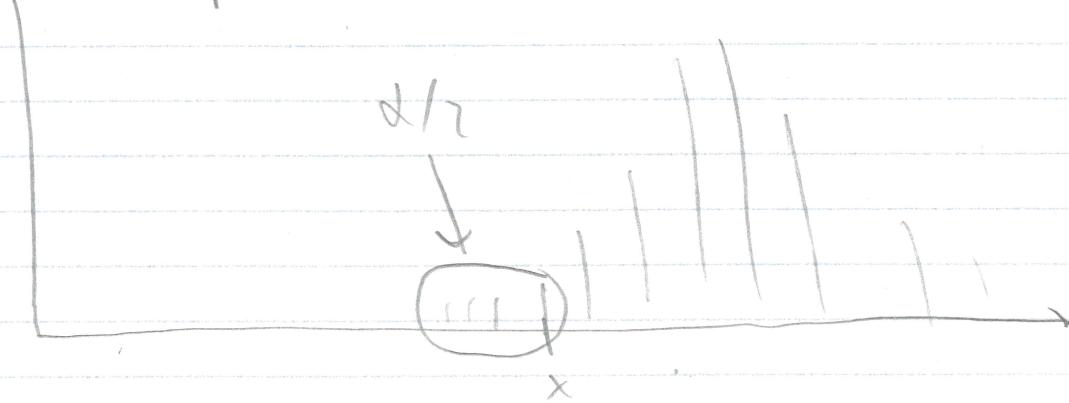
$$\frac{x}{2} = \text{pbinom}(x, \text{size}=n, \text{prob}=p_1)$$

$$\frac{x}{2} = 1 - \text{pbinom}(x-1, \text{size}=n, \text{prob}=p_2)$$

under  $p_2$



under  $p_1$



Binomial R code

## • Poisson Estimation

$X \sim \text{Poi}(\mu)$  where  $\mu = \lambda T$

Goal: Estimate  $\lambda$

$\lambda$  = Incidence rate per unit time

$T$  = time

Ex.) Leukemia rate in Woburn, MA. 12,000 residents, 10 years. 12 children got Leukemia.

• A common unit of time in Biostatistics is a person-year

↳ 1 person being followed for 1 year

↳ Normalizes by # people in a study

↑ Eg. 12 cases out of 20 is a lot more than  
12 out of 20,000

Ex.) Woburn study had 12,000 people • 10 years  
= 120,000 person-years.

$$\hat{\lambda} = \frac{x}{T}, \quad E(\hat{\lambda}) = \frac{E(x)}{T} = \frac{\lambda T}{T} = \lambda$$

Unbiased

$$\text{Woburn: } \hat{\lambda} = \frac{12 \text{ cases}}{120,000 \text{ person years}} = 0.001 \frac{\text{cases}}{\text{person year}}$$

- Conc. rates low so need to multiply by 100,000

$$\frac{0.0001 \text{ cases}}{\text{person-year}} = 10 \frac{\text{cases}}{100,000 \text{ person-years}}$$

- Interval size strategy as exact Binomial case

$$p_1 \text{ st } \Pr(X \geq x | p_1) = \alpha/2$$

$$p_2 \text{ s.t. } \Pr(X \leq x | p_2) = 2/3$$

Interval  $(\mu_1, \mu_2)$

## CDF of Poisson

- Interval for  $\lambda$  is  $\left(\frac{\mu_1}{T}, \frac{\mu_2}{T}\right)$

- Use `pairson.test()` in R

## • one-sided CI (UK to ship)

↑ only interested in lower or upper bounds

Ex: compare cancer treatment survival rate against baseline of 30%. Want lower bound on treatment effect to see if better than baseline

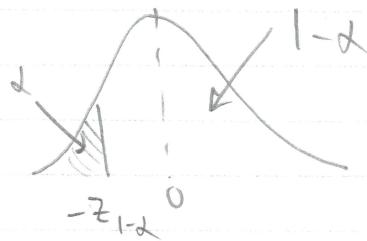
• Fld  $p_1$  s.t.  $P_r(p > p_1) = 1-\alpha$

↑ random lower bound

• Or, fnd  $p_2$  s.t.  $P_r(p < p_2) = 1-\alpha$

↑ Ex: Upper bound on incidence rate of some treatment group

• Since  $\frac{\hat{p} - p}{\sqrt{\hat{p}(1-\hat{p})/n}} \sim N(0,1)$



$$P_r\left(-z_{1-\alpha} < \frac{\hat{p} - p}{\sqrt{\hat{p}(1-\hat{p})/n}}\right) = 1-\alpha$$

$$\Rightarrow P_r\left(z_{1-\alpha} > \frac{p - \hat{p}}{\sqrt{\hat{p}(1-\hat{p})/n}}\right) = 1-\alpha$$

$$\Rightarrow P_r\left(\hat{p} + z_{1-\alpha} \sqrt{\hat{p}(1-\hat{p})/n} > p\right) = 1-\alpha$$

$$\Rightarrow p_2 = \hat{p} + z_{1-\alpha} \sqrt{\hat{p}(1-\hat{p})/n}$$

$$\text{Similarly } p_1 = \hat{p} - z_{1-\alpha} \sqrt{\hat{p}(1-\hat{p})/n}$$

Ex.) 40 out of 100 patients survive a new cancer treatment. What is lower bound on  $\Pr(\text{survive})$ ? (Upper 1-sided 95% CI)

$$\hat{p} = \frac{40}{100} = 0.4$$

$$z_{1-\alpha} = \text{qnorm}(0.95) = 1.645$$

$$\alpha = 0.05$$

$$\hat{p} - z_{1-\alpha} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

$$= 0.4 - 1.645 \sqrt{\frac{0.4 \cdot 0.6}{100}}$$

$$\approx 0.319$$