

## December 1, Lecture 21

~~Recall if  $X_1, \dots, X_n$  are iid with~~

Recall if  $X_1, \dots, X_n$  iid with mean  $\mu$  & standard error  $\sigma$ , and  $n$  large...

$$(II) \frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}}$$

$$(III) \bar{X} \stackrel{d}{\approx} N(\mu, (\frac{\sigma}{\sqrt{n}})^2)$$

$$(IV) T \stackrel{d}{\approx} N(n\mu, (\sqrt{n} \sigma)^2)$$

Shipments are late 2% of the time. In 10,000 shipments, what is the probability more than 3% are late?

\*  $X_1, \dots, X_{10000} \stackrel{iid}{\sim} \text{Bern}(p=2\%)$

→ Assumption: assuming each shipment iid Bern( $p=2\%$ ).

→ Assume it's right.

\*  $P(\bar{X} \geq 3\%) = P(\hat{p} \geq 0.3) \rightarrow \text{next page} \dots$

\*  $\bar{X} \stackrel{\text{by III}}{\approx} N(\mu, (\frac{\sigma}{\sqrt{n}})^2) \rightarrow \text{If } X_1, \dots, X_n, \text{ what is } \mu \text{ \& } \sigma?$

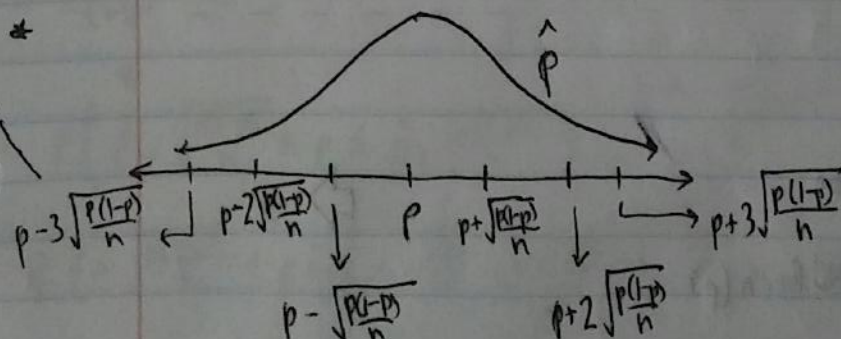
$$* \hat{p} \sim N(p, (\sqrt{\frac{p(1-p)}{n}})^2)$$

$\hat{p} = \bar{X}$  in Bernoulli case

$$* \mu = p$$

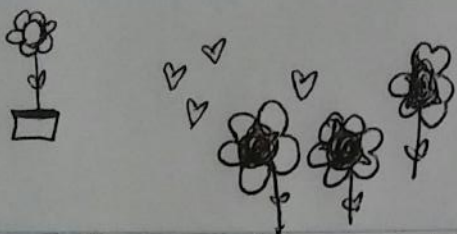
$$* \sigma = \sqrt{p(1-p)}$$

$$* \hat{p} := \bar{X} = \frac{\sum \#1's}{n} \quad \text{"sample proportion"}$$



→ Next





$$P(\bar{X} \geq 3\%) = P(\hat{p} \geq 0.03) = P\left(\frac{\hat{p} - .02}{\sqrt{\frac{.02(1-.02)}{10000}}} \geq \frac{.03 - .02}{\sqrt{\frac{.02(1-.02)}{10000}}}\right)$$

$$P(Z \geq 7.14) \approx 0$$

Implication - Never going to have more than 3% late

• Mushrooms  $\rightarrow$  11 yes, 23 total subjects

$$n = 23$$

$$\hat{p} = 11/23 = 0.48 \text{ ("data")}$$

$[\hat{p} \text{ came from } p]$  - little  $\hat{p}$  with no hat is  $\Rightarrow$  true expectation of someone like mushrooms or not.

\* "p" is the "true" population parameter. \*

\* We don't know p, too many people to ask. Need to ask everyone, however with a population of 7.5 billion, too difficult. Unknowable \*

[Our goal is to know something about p.  
Can we say anything we know about ~~our~~ p.]

Statistics has started, probability ended.

We have data, but don't know parameter.

Trying to infer something from data about parameter.

\* Statistical Inference: Infer the population parameter using a statistic of the data. Here,  $\hat{p}$ .

population, size  $N \propto \infty$

sample size  
of  $n < \infty$

$X_1, \dots, X_n \stackrel{iid}{\sim} \text{Bern}(p)$

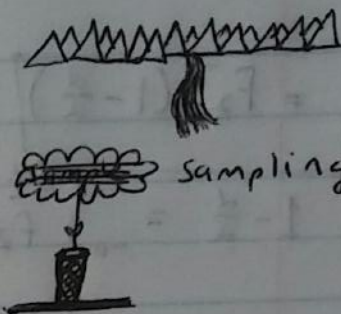
$X_1, \dots, X_n \stackrel{iid}{\sim} \text{Bern}(p)$





\*What constitutes a "good" sample? \*

A representative sample. ie completely random sampling.



## \*Goals of Inference

- ① Give me the best guess of  $p$ .
- ② Give me a ~~reasonable~~ reasonable interval of values for  $p$ .
- ③ Test theories about  $p$ .

### ① point estimation

$$p \approx \hat{p}$$

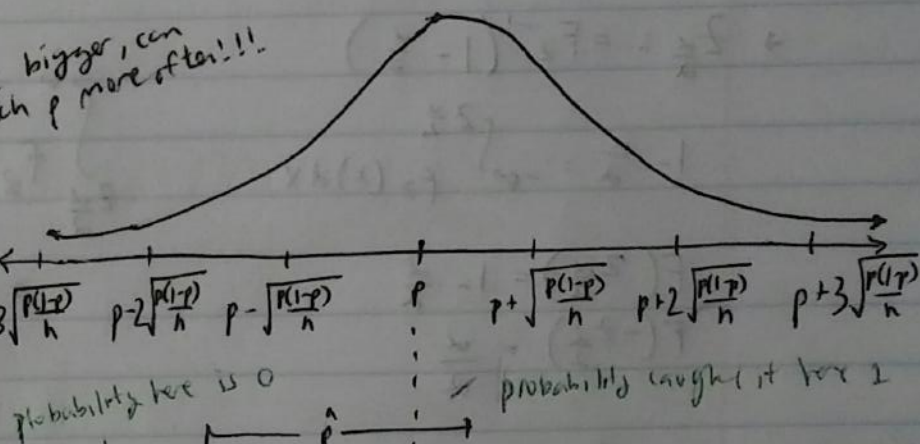
### ② Confidence Intervals

- If make bigger, can catch  $p$  more often!!!



$$\left[ \hat{p} \pm z_{\frac{\alpha}{2}} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \right]$$

have to amplify, make bigger, to catch  $p$  more.



Did this interval catch true  $p$ ? Is  $p \in \left[ \hat{p} \pm \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \right]$ ?

→ Yes, because drew it

that way.

If do this many many times,

how often will catch the truth? - What's probability that catch it?

$$\left[ \hat{p} \pm \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \right] := \left[ \hat{p} - \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}, \hat{p} + \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \right]$$

$$P\left(p \in \left[ \hat{p} \pm \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \right]\right) \quad \hat{p}, \hat{p}, p, P()$$

(big)

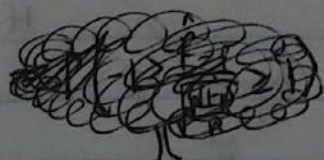
$$P\left(\hat{p} - \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \leq p \leq \hat{p} + \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}\right)$$

$$P\left(-\sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \leq p - \hat{p} \leq \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}\right)$$

$$= P\left(-1 \leq \frac{p - \hat{p}}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}} \leq 1\right)$$

CLT

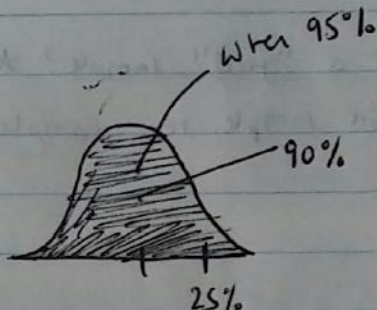
$$P(-1 \leq -Z \leq 1) = P(Z \in [-1, 1]) = 68\%$$





$$z_{\frac{\alpha}{2}} = F_z^{-1}\left(1 - \frac{\alpha}{2}\right)$$

$$1 - \frac{\alpha}{2} = \int_{-\infty}^{z_{\frac{\alpha}{2}}} f_z(x) dx$$



$$\alpha = 10\% \Rightarrow \frac{\alpha}{2} = 5\%$$

$$1 - \frac{\alpha}{2} = 95\%$$

If do calculation again...

$$* \left[ \hat{p} \pm z_{\frac{\alpha}{2}} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \right] \dots \rightarrow P\left(z \in \left[-z_{\frac{\alpha}{2}}, z_{\frac{\alpha}{2}}\right]\right) \leftarrow F\left(z_{\frac{\alpha}{2}}\right) - F\left(-z_{\frac{\alpha}{2}}\right) \dots \left(1 - \frac{\alpha}{2}\right) - \frac{\alpha}{2} = 1 - \alpha$$

$$\Rightarrow z_{\frac{\alpha}{2}} = F_z^{-1}\left(1 - \frac{\alpha}{2}\right)$$

$$1 - \frac{\alpha}{2} = \int_{-\infty}^{z_{\frac{\alpha}{2}}} f_z(x) dx$$

$$F\left(z_{\frac{\alpha}{2}}\right) = 1 - \frac{\alpha}{2}$$

$$F\left(-z_{\frac{\alpha}{2}}\right) = \frac{\alpha}{2}$$

$$\int_{z_{\frac{\alpha}{2}}}^{\infty} f_z(x) dx = 1 - \left(1 - \frac{\alpha}{2}\right) = \frac{\alpha}{2} = \int_{-\infty}^{-z_{\frac{\alpha}{2}}} f_z(x) dx$$

\* Mushroom example

$$\hat{p} = 0.48$$

~~0.48 ± 2 for 95% coverage~~

$$\left[ 0.48 \pm 2 \sqrt{\frac{p(1-p)}{28}} \right] \text{ for 95\% coverage.}$$

absurd, don't know p. What we are trying to look for anyway.

$$* \left[ \hat{p} \pm z_{\frac{\alpha}{2}} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \right] \approx \left[ \hat{p} \pm z_{\frac{\alpha}{2}} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \right] \Rightarrow \left[ 0.48 \pm 2 \sqrt{\frac{0.48(1-0.48)}{23}} \right]$$

change to

If the p is not  $p \approx 0$  or  $p \approx 1$ .

debated for 100 years.

\* Definition: A  $1 - \alpha$  size confidence interval for population proportion p

$$CI_{1-\alpha, p} = \left[ \hat{p} \pm z_{\frac{\alpha}{2}} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \right]$$

$$* [0.48 \pm 2 \sqrt{\frac{0.48(1-0.48)}{23}}] \Rightarrow [0.272, 0.688]$$

for 95%

coverage

\* Can you say this?  $P(p \in [0.272, 0.688]) \neq 95\%$

No, the real  $p$  is a single #. This is

range of #s & is not random.

So no, can't say this.