

21

keyword implied average

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

model $X_1 + X_2 + \dots + X_{10,000} \stackrel{\text{iid}}{\sim} \text{Bernoulli}(2\%)$

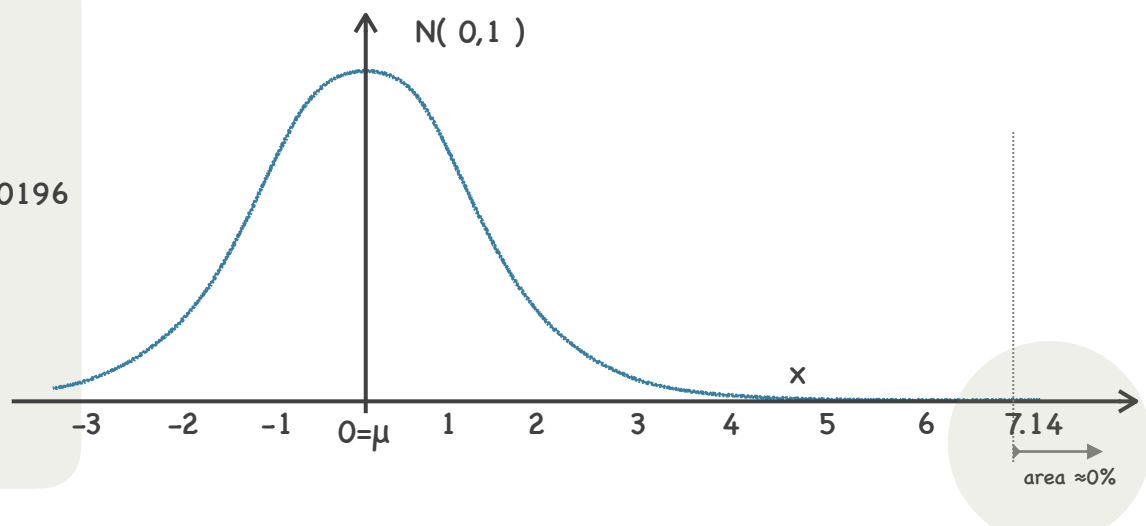
ideally we must assume iid, otherwise we can't do anything

$\mu = p = 0.02$
 $\sigma^2 = p(1-p) = 0.02 \cdot 0.98 = 0.0196$
 $\Rightarrow \sigma = \sqrt{0.0196} = 0.14$
S.E. = $\frac{\sigma}{\sqrt{n}} = \frac{0.14}{\sqrt{10,000}} = .0014$

$\bar{X} \approx \text{Normal}(\mu, (\frac{\sigma}{\sqrt{n}})^2) = \text{Normal}(.02, .0014^2)$

probability statement $P(\bar{X} > 3\%) = P(\frac{\bar{X} - .02}{.0014} > \frac{.03 - .02}{.0014}) \approx P(Z > 7.14) \approx 0$

standardization by CLT
subtract the mean and divide by S.E.



P-Hat

let's define a new r.v. called **P-Hat** 'Sample Proportion'

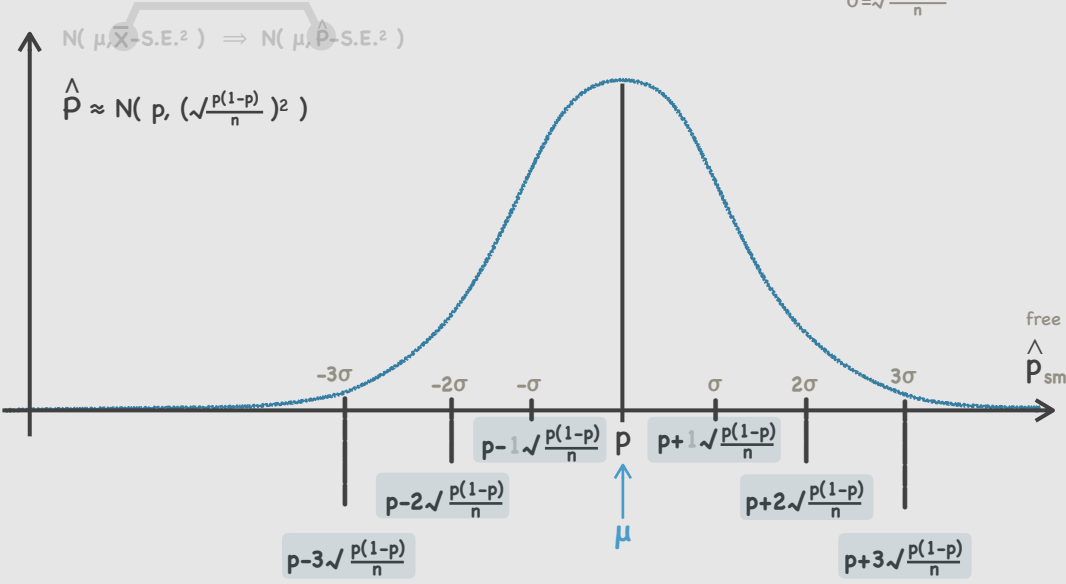
upper-case $\hat{p} = \bar{X}$ 'P-hat' - Sample proportion r.v.
lower-case $\hat{p} := \bar{x} = \frac{\sum x_i}{n} = \frac{\# \text{ of 1s}}{n}$ in Bernoulli

little p-hat is locked between 0 and 1 subset of all averages

$\bar{X} \approx \text{Normal}(\mu, (\frac{\sigma}{\sqrt{n}})^2)$
 $\hat{p} \approx \text{Normal}(\mu, (\frac{\sigma}{\sqrt{n}})^2)$

Bernoulli $\hat{p} \approx \text{Normal}(p, (\sqrt{\frac{p(1-p)}{n}})^2)$

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



P-Hat is a 'normal distribution'

probability statement $P(\hat{p} > 3\%) = P(\frac{\hat{p} - .02}{.0014} > \frac{.03 - .02}{.0014}) \approx P(Z > 7.14) \approx 0$

standardization by CLT
subtract the mean and divide by S.E.
 $\mu = p = 0.02$
 $\sigma = \sqrt{0.02(1-0.02)} = .0014 \approx \text{S.E.}$

thus, 3% late will never happen

Who likes mushrooms? $\hat{p} = \frac{\# \text{ of students who like mushrooms}}{\text{total \# of sampled students}} = \frac{11}{23} = 0.48 = 48\%$

What is 'p'? 'p' is a true expectation of someone liking mushrooms, $p = \mu$

If there is 7.5 billion people and we were to sample every single one of them, then we could find our 'p'. Thus 'p' is the true population parameter but since it is neither practical nor realistic to sample all 7.5 billion people, 'p' is unknowable.

However, our goal is to know something about 'p'.

Statistics

Can we use our classroom sampling to know something about 'p'?

Previously we were given r.v. models with all the parameter values. We were able to calculate data based on those knowable quantities of the parameters. Now we are facing the inverse of the problem. We have data but we do not know the parameters. We are trying to infer the parameters from the acquired data.

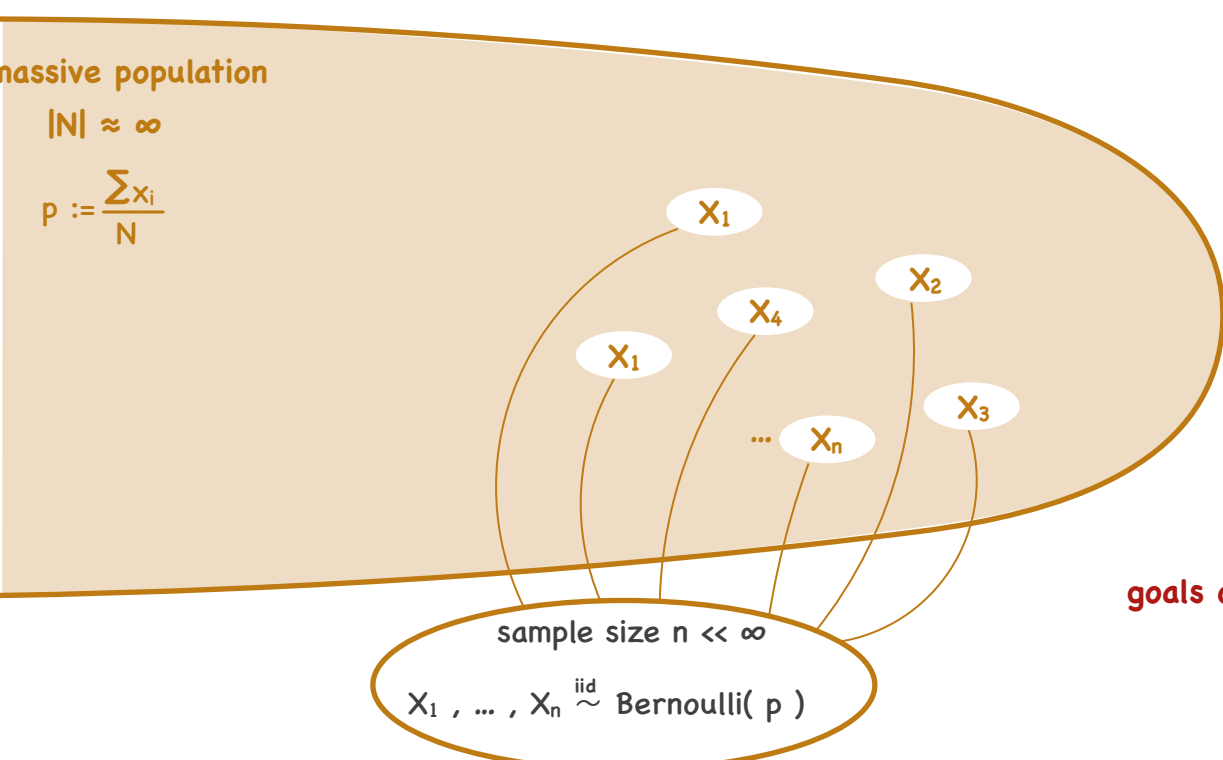
no given parameters \Rightarrow infer the parameters from data

Statistical Inference: infer population parameter using the statistics of the data

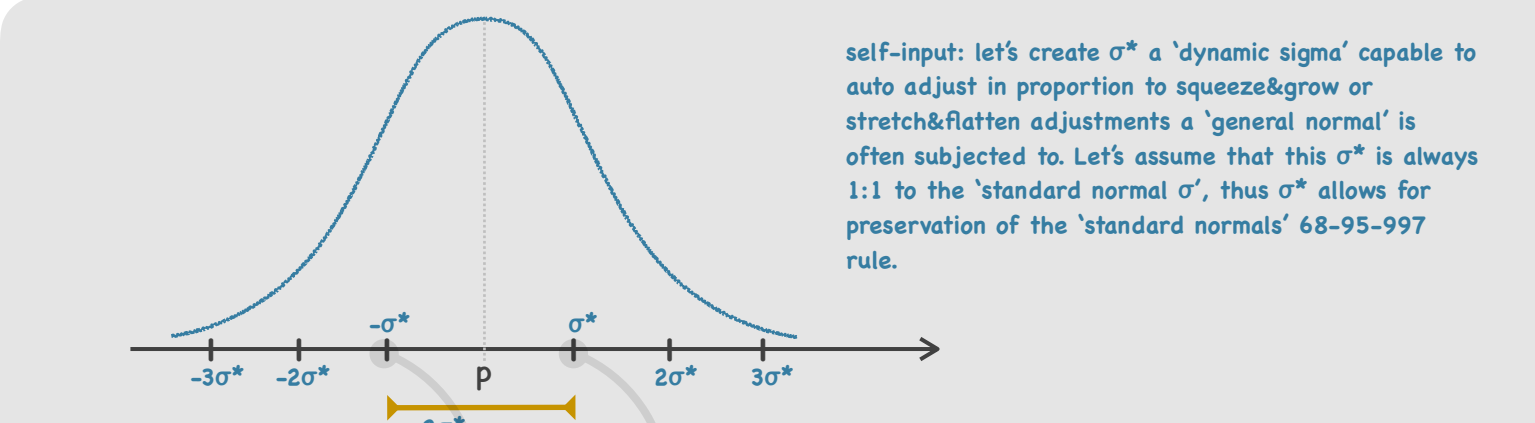
In order to know something about the truth of 'p' we can collect a 'finite sample' or 'small sample', and then use it. What constitutes a good sample? Sample must be 'representative' which means it preserves iid propensity. How? Simple random sample. All males? All college students? No... it must be completely random. (attempt at the encapsulation of the entire gamut of diversity)

goals of inference:

1. give me the best guess of p - 'point estimation' (estimate p as a single point)
2. give me a reasonable interval of values for p - 'interval construction' (estimate a range of ps which makes sense)
3. let me test theories about p (test theories about what p is)



Interval Construction aka 'Confidence Intervals'



lets create an interval

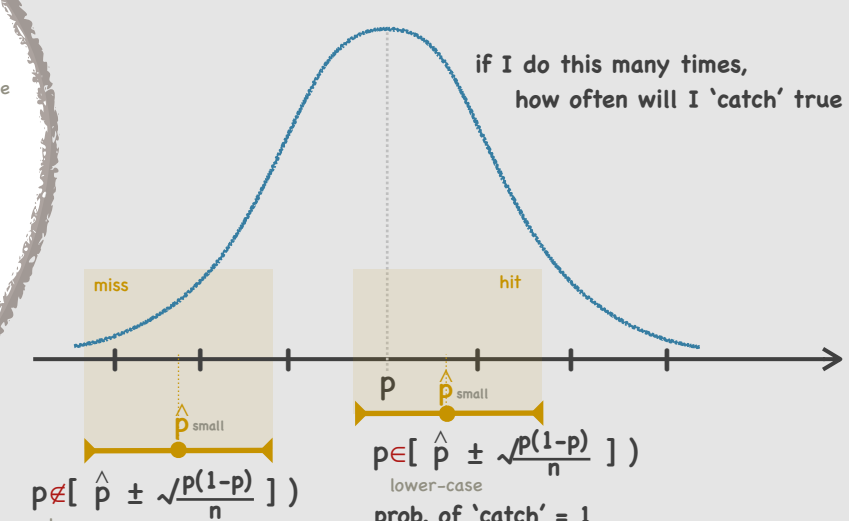
$$[\hat{p} \pm \sqrt{\frac{p(1-p)}{n}}] := [\hat{p} - \sqrt{\frac{p(1-p)}{n}}, \hat{p} + \sqrt{\frac{p(1-p)}{n}}]$$

analogous to saying $[5 \pm 1] = [4, 6]$ $2\sigma^*$ range

What is the probability that 'p' is in the range of \hat{p} . Did this interval capture the true 'p'?

upper-case $P(p \in [\hat{p} \pm \sqrt{\frac{p(1-p)}{n}}])$
lower-case $p \in [\hat{p} \pm \sqrt{\frac{p(1-p)}{n}}]$
Expectation 'True-p'

what is the prob. that I'll 'catch' the 'true-p' using this interval?
we are interested in the long run probability



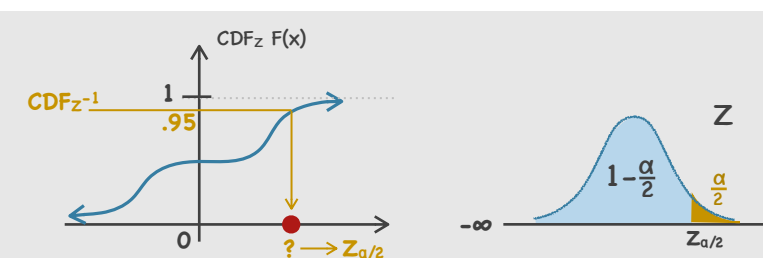
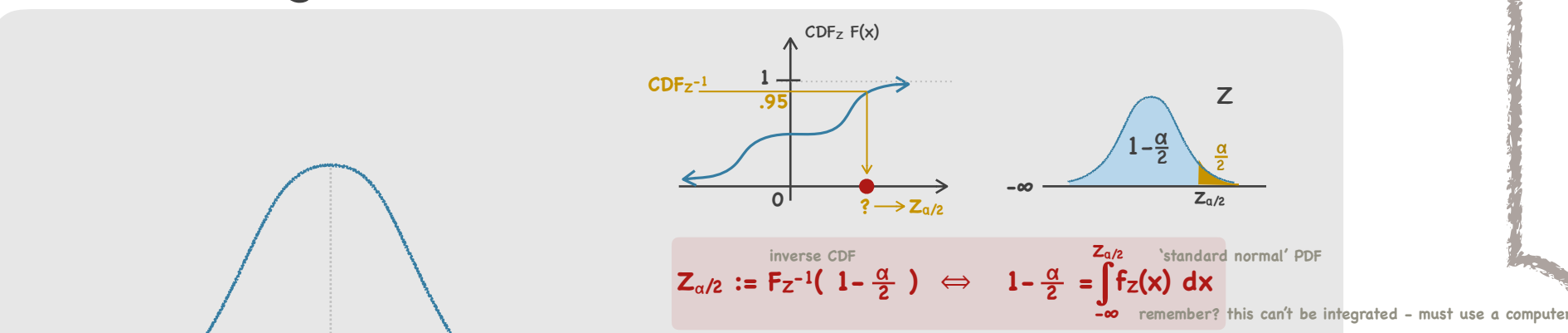
if I do this many times, how often will I 'catch' true 'p'

upper-case $P(\hat{p} - \sqrt{\frac{p(1-p)}{n}} \leq p \leq \hat{p} + \sqrt{\frac{p(1-p)}{n}})$
lower-case $p \in [\hat{p} \pm \sqrt{\frac{p(1-p)}{n}}]$
 $= P(-1 \leq \frac{\hat{p} - p}{\sqrt{\frac{p(1-p)}{n}}} \leq 1)$
 $= P(-1 \leq Z \leq 1)$
 $= P(1 \geq Z \geq -1)$
 $= P(-1 \leq Z \leq 1)$
 $= P(Z \in [-1, 1]) = .68$
remember the quantiles '68-95-997' rule? $[-\sigma^*, \sigma^*]$

thus by creating this interval we will 'catch' the 'true-p' 68% of the time (utopia)

lets create a bigger 'paddle'

Larger Interval Construction



what does this mean?
 $\alpha = 10\% \Rightarrow \frac{\alpha}{2} = 5\% \Rightarrow 1 - \frac{\alpha}{2} = 95\%$
 $\alpha = 5\% \Rightarrow \frac{\alpha}{2} = 2.5\% \Rightarrow 1 - \frac{\alpha}{2} = 97.5\%$

analogous to saying $[5 \pm 2] = [3, 7]$ $4\sigma^*$ range

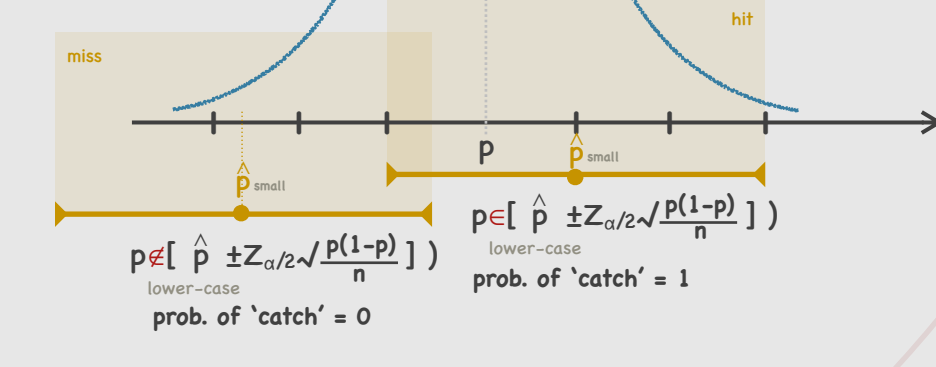
$$[\hat{p} \pm Z_{\alpha/2} \sqrt{\frac{p(1-p)}{n}}] := [\hat{p} - Z_{\alpha/2} \sqrt{\frac{p(1-p)}{n}}, \hat{p} + Z_{\alpha/2} \sqrt{\frac{p(1-p)}{n}}]$$

What is the probability that 'p' is in the range of \hat{p} . Did this larger interval capture the true 'p'?

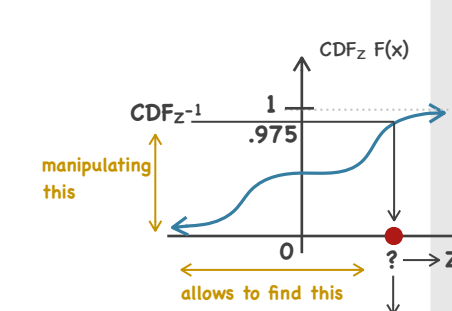
upper-case $P(p \in [\hat{p} \pm Z_{\alpha/2} \sqrt{\frac{p(1-p)}{n}}])$
lower-case $p \in [\hat{p} \pm Z_{\alpha/2} \sqrt{\frac{p(1-p)}{n}}]$
Expectation 'p'

if I do this many times, how often will I 'catch' the 'true-p'

upper-case $P(\hat{p} - Z_{\alpha/2} \sqrt{\frac{p(1-p)}{n}} \leq p \leq \hat{p} + Z_{\alpha/2} \sqrt{\frac{p(1-p)}{n}})$
lower-case $p \in [\hat{p} \pm Z_{\alpha/2} \sqrt{\frac{p(1-p)}{n}}]$
 $= P(-Z_{\alpha/2} \leq \frac{\hat{p} - p}{\sqrt{\frac{p(1-p)}{n}}} \leq Z_{\alpha/2})$
CLT $= P(-\frac{\alpha}{2} \leq Z \leq \frac{\alpha}{2})$
 $= P(\frac{\alpha}{2} \geq Z \geq -\frac{\alpha}{2})$
 $= P(-\frac{\alpha}{2} \leq Z \leq \frac{\alpha}{2})$
 $= P(Z \in [-\frac{\alpha}{2}, \frac{\alpha}{2}]) = .95$
remember the quantiles '68-95-997' rule? $[-2\sigma^*, 2\sigma^*]$



thus by creating this bigger interval we hope to 'catch' the 'true-p' 95% of the time



$F(Z_{\alpha/2}) = 1 - \frac{\alpha}{2}$
 $F(-Z_{\alpha/2}) = \frac{\alpha}{2}$
 $F(Z_{2.5}) = 97.5\%$
 $F_{Z^{-1}}(.975) = Z_{2.5}$
now we can operate on a limited budget especially if 80% is informative enough. Let's play 'pong' with a paddle large enough to win. Meaning we can now make the 'paddle' size according to what we expect to get.

Point Estimation

How to get the best guess of p? $\hat{p} := \frac{\sum x_i}{n} = \frac{\# \text{ of 1s}}{n}$ sample proportion

Where does \hat{p} come from? \hat{p} is a realization from \hat{p}

$$p \approx \hat{p} = \frac{11}{23} = 0.48 = 48\% \text{ (ppl like mushroom)}$$

the Classic Method

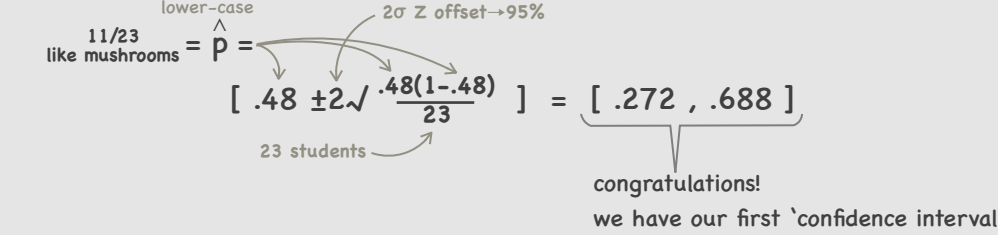
and the solution is ... to use a smart approximation ... debated for 100 years (and still is)

$$[\hat{p} \pm Z_{\alpha/2} \sqrt{\frac{p(1-p)}{n}}] \approx [\hat{p} \pm Z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}]$$

as long as $p \neq 0$ and $p \neq 1$

Def: a $1-\alpha$ sized 'confidence interval' for population proportion p is:

$$CI_{p, 1-\alpha} := [\hat{p} \pm Z_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}]$$



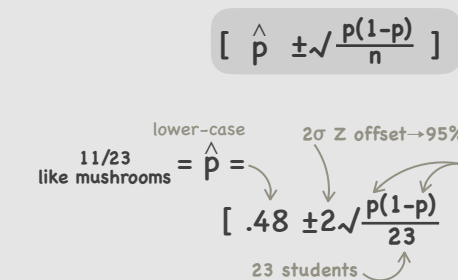
What does this interval mean? Can we be hopeful that the 'true-p' lies somewhere within this interval?
Can we say this:

$$P(p \in [.272, .688])$$

Unfortunately we can NOT say this! b/c this might be completely false. Our 48% comes from a single sample. This is like trying to calculate a mean of a single try at a roulette table.

Let's catch that p

Great! Let's make a custom 'paddle' that will 'catch' the 'true-p' 95% of the time



paradox - in order to find the 'true-p' 95% of the time, we need the 'true-p'. But if we did know the 'true-p' to start with, we wouldn't be looking for it.