

# The Normal Approximation $\hat{=}$ Binomial Proportion Confidence Interval

## Central Limit Theorem

Suppose  $\{X_i\}_{i=1}^n$  are any collection of iid random variables with mean  $\mu = \mathbb{E}X_i < \infty$  and standard deviation  $\sigma = \text{sd}(X_i) < \infty$ .

Then the sum  $S_n = \sum_{i=1}^n X_i$  and mean  $\frac{S_n}{n}$  converge in distribution to the normal distribution,

$$\frac{S_n - n\mu}{\sigma\sqrt{n}} \xrightarrow{d} \mathcal{N}(0,1) \quad \text{as } n \rightarrow \infty$$

$$\frac{\frac{S_n}{n} - \mu}{\sigma/\sqrt{n}} \xrightarrow{d} \mathcal{N}(0,1) \quad \text{as } n \rightarrow \infty$$

Convergence in distribution  $\xrightarrow{d}$  means

$$P(a \leq \frac{S_n - n\mu}{\sigma\sqrt{n}} \leq b) \rightarrow P(a \leq Z \leq b) \quad \text{as } n \rightarrow \infty,$$

Tons of quantities in sports are the sum or mean of iid random variables! So the normal approximation comes in handy!

# Estimating the quality of a free throw shooter

Shaq shoots  $n$  free throws, whose results

are given by  $\{X_i\}_{i=1}^n$ ,  $X_i = \begin{cases} 1 & \text{if } i\text{th free} \\ & \text{throw made} \\ 0 & \text{if missed} \end{cases}$

$S_n = \sum_{i=1}^n X_i$  is his # made free throws.

We model  $S_n$  by

$S_n \sim \text{Binomial}(n, p) =$  # successes in  $n$  trials  
where each trial is  
independent Bernoulli( $p$ )  
( $X_i \stackrel{\text{iid}}{\sim} \text{Ber}(p)$ )

We want to estimate Shaq's  $p$  from the data  $\{X_i\}_{i=1}^n$ .

Our "best guess" of  $p$  is  $\hat{p} = \frac{S_n}{n} = \frac{1}{n} \sum_{i=1}^n X_i$ .

In fact, this is the **MLE** (see the Lab).

But how confident should we be in this guess?

Let  $\mu = \mathbb{E}X_i = p$ ,  $\sigma^2 = \text{VAR}(X_i) = \sqrt{p(1-p)}$ .

By CLT,  $\frac{\hat{p} - \mu}{\sigma/\sqrt{n}} \xrightarrow{d} \mathcal{N}(0,1)$  as  $n \rightarrow \infty$ .

So  $P(-2 \leq \frac{\hat{p} - \mu}{\sigma/\sqrt{n}} \leq 2) = P(-2 \leq \frac{p - \hat{p}}{\sqrt{p(1-p)}/\sqrt{n}} \leq 2) \approx 0.95$

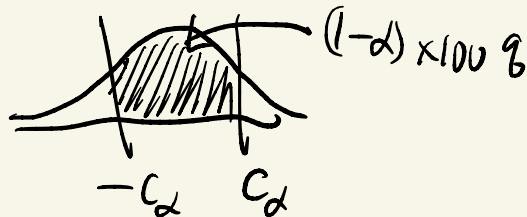
So  $p \in \hat{p} \pm 2 \sqrt{p(1-p)/n}$  w.p.  $\approx 0.95$

$p$  is unknown, so use  $\hat{p}$  in place of  $p$

Wald CI:  $\hat{p} \pm 2 \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$

More generally:  $\hat{p} \pm C_\alpha \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$

is a  $(1-\alpha) \times 100\%$  CI where  $C_\alpha$  is the quantile of  $\mathcal{N}(0,1)$  taking  $(1-\alpha) \times 100\%$  of the area



You buy an M&M bag with 56 M&M's and 14 blue ones. Supposing the color of each M&M is randomly drawn from some iid distribution, what is a 95% confidence interval of the probability the company makes an M&M blue?

$$\hat{p} = \frac{14}{56} = \frac{1}{4}$$

$$\hat{p} \pm 2 \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} = \frac{1}{4} \pm 2 \sqrt{\frac{\frac{1}{4} \cdot \frac{3}{4}}{n}} = \frac{1}{4} \pm \left(\frac{\sqrt{3}}{2}\right) \cdot \frac{1}{\sqrt{n}} \approx .85$$

$$n=56 \Rightarrow CI \approx \frac{1}{4} \pm \frac{1}{2} \quad \text{pretty wide!}$$

$$n=400 \Rightarrow CI \approx \frac{1}{4} \pm .0425$$

# The Meaning of a confidence interval

The confidence interval of a Binomial proportion  $p$  is

$$\text{Wald CI} = \hat{p} \pm 2 \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

$$\text{i.e. } P(p \in \text{CI}) \approx 0.95$$

but what does this probability  $P$  actually mean?

Under this model,  $S_n = \sum_{i=1}^n X_i \sim \text{Binomial}(n, p)$

$X_i \stackrel{\text{iid}}{\sim} \text{Bernoulli}(p)$

(Frequentist)  
↑

Under this model,  $p$  is an unknown fixed constant.

$p$  either lies in the CI or doesn't lie in it, so in actuality the probability it lies in the CI is either 0 or 1. So, what does  $P$  mean?

The probability  $P$  refers to randomness in the CI:

The CI itself is a random variable because  $\hat{p}$  is a random variable because  $\{X_i\}_{i=1}^n$  are random variables.

If we replicated the experiment 100 times, in each replication  $p$  remains the same but the data  $\{X_i\}$  and the CI's change by randomness.

the CI is expected to contain the true <sup>on average</sup> unobserved  $p$  in 95 of these 100 replications

## Coverage

The Wald CI is based on 2 approximations:

that  $P\left(-2 \leq \frac{\hat{p} - p}{\sqrt{p(1-p)/n}} \leq 2\right) \approx 0.95$  by CLT

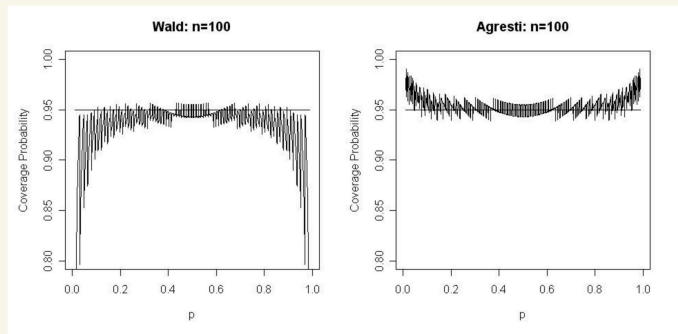
and that we can plug in  $\hat{p}$  for  $p$  in  $\sqrt{p(1-p)}$

so that  $P\left(-2 \leq \frac{\hat{p} - p}{\sqrt{\hat{p}(1-\hat{p})/n}} \leq 2\right) \approx 0.95.$

The **actual** coverage Probability  $IP(p \in CI)$  of the 95% Wald interval can be quite far below the **nominal** coverage of 95% as shown by simulations and computations (Brown, Cai, Dasgupta 2001).

If  $n$  is several hundred or thousand the Wald interval is tolerably accurate. Otherwise we need to adjust the CI.

Agresti and Coull (1998) recommend introducing **2 artificial successes and failures** into the data before computing  $\hat{p}$  and  $n$ , which is better than the Wald interval.



→ oscillations from the discreteness of Binomial

Agresti-Coull  $CI = \hat{p}' \pm 2 \sqrt{\frac{\hat{p}'(1-\hat{p}')}{n+4}}$  where  $\hat{p}' = \frac{S_n+2}{n+4}$ .