# Data Science Example Social Ratings

## Last time

Model social ratings (thumbs up/down) as $n$ $0,1$ (Bernoulli) variables

Rating $\bar{X} = \hat{p} = \frac{k}{n}$ "well" approximated by <u>normal</u> distribution

Use Normal's <u>Confidence Intervals</u> to determine likely/unlikely values of $p$ (item's true rating) given $\hat{p}, n$.

$$\hat{P}_L = P - 1.96\sqrt{\tfrac{1}{n}P(1-P)} \qquad \text{mean} - 1.96\,(\text{stdv})$$

$$\hat{P}_R = P + 1.96\sqrt{\tfrac{1}{n}P(1-P)} \qquad \text{mean} + 1.96\,(\text{stdv})$$

$$\rightarrow 1.96 \Rightarrow 95\% \text{ C.I.}$$

**\* Lower Confidence Bound (LCB) sort**

- Sort items using $\hat{P}_L$ (a "worst-case" estimate of $p$)

- Statistically well-motivated way to combine our sort objectives $\Rightarrow$ Rank items incorporating uncertainty

## Remaining Limitation

$$\hat{P}_L = P - 1.96\sqrt{\tfrac{1}{n}P(1-P)} \quad \text{depends on } p, \; p \text{ unknown!}$$

How to compute $\hat{P}_L$?

Here's two solutions.

    1. Use sample statistics $\rightarrow$ replace $p$ w/ $\hat{p}$ in $\hat{P}_L$.

$$\hat{P}_L = \hat{p} - 1.96\sqrt{S_x^2}$$

         &uarr; sample mean = $\hat{p}$      &uarr; sample variance = $S_x^2$

    **\* Sometimes called "Wald Approximation" Can be OK to use but not always accurate.**

    2. <u>Wilson Score</u> - Let's study this for some nice insights. $\rightarrow$

$\boxed{1}$

## Wilson Score (w.s.)

Let $\pm z = \dfrac{\hat{P} - P}{\sqrt{\dfrac{P(1-P)}{n}}}$      95% C.I on $z$ is $[-1.96, 1.96]$

w.s. → solve this for $P$:

$$\frac{\hat{P} + \dfrac{z^2}{2n} \pm z\sqrt{\dfrac{\hat{P}(1-\hat{P}) + \dfrac{z^2}{4n}}{n}}}{1 + \dfrac{z^2}{n}} = P \quad \leftarrow \text{Get } P_L, P_R \text{ by plugging in } \hat{P}, n, z$$

---

**That's the answer but let's dig deeper:**

Let's rewrite this to understand it better.

Here $P$ is of the form $A \pm B$, let's focus on $A$ (stuff left of $\pm$):

$$\frac{\hat{P} + \dfrac{z^2}{2n}}{1 + \dfrac{z^2}{n}}$$

- recall $\hat{P} = \dfrac{k}{n}$   $k$ t.u.'s of $n$ ratings
  → plug in
- Also, Lets plug in $z = 2 \approx z_L = 1.96$ $^{\text{close enough!}}$

$$\approx \frac{\dfrac{k}{n} + \dfrac{4}{2n}}{1 + \dfrac{4}{n}} = \frac{\dfrac{k+2}{n}}{\dfrac{n+4}{n}} = \frac{k+2}{n+4}$$

⇒ Wilson score is a __smoothed__ approximation!
  add 2 successes and 2 failures   $k \to k+2$
  $\underbrace{\phantom{2 \text{ successes}}}_{t.u.}$  $\underbrace{\phantom{2 \text{ failures}}}_{t.d.}$   $n \to n+4$

⇒ This idea of "smoothing" low count data is very common. Can appear __ad hoc__ but in many situations is statistically well principled (of course, here we only looked at term to left of $\pm$).