# Lecture 13: Confidence Intervals (II)

Ye Tian

Department of Statistics, Columbia University
Calculus-based Introduction to Statistics (S1201)

Aug 1, 2022

# Recap: Confidence intervals

**<u>Definition</u>**: Suppose we have samples $X_1, \ldots, X_n$. If we can construct an random interval $[l(X_1, \ldots, X_n), u(X_1, \ldots, X_n)]$ which satisfies

$$\mathbb{P}(l(X_1, \ldots, X_n) \leq \theta \leq u(X_1, \ldots, X_n)) = 1 - \alpha,$$

then it can be called as a $(100(1 - \alpha))\%$ **confidence interval (CI)**. $(100(1 - \alpha))\%$ is the **confidence level**.

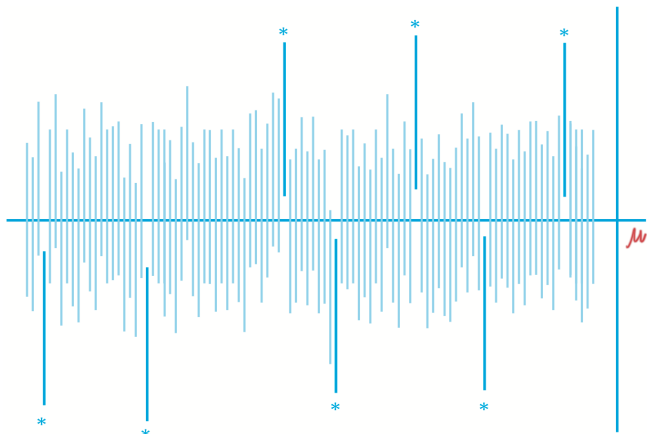**<u>Example</u>**: $X_1, \ldots, X_n \overset{i.i.d.}{\sim} N(\mu, 1)$

$$\mathbb{P}\left( \bar{X} - \frac{z_{0.025}}{\sqrt{n}} \leq \mu \leq \bar{X} + \frac{z_{0.025}}{\sqrt{n}} \right) = 95\%.$$

Thus the 95% CI is $\left[ \bar{X} - \frac{z_{0.025}}{\sqrt{n}}, \bar{X} + \frac{z_{0.025}}{\sqrt{n}} \right]$.

# How to understand this 95% coverage probability

$$\mathbb{P}\left(\bar{X} - \frac{z_{0.025}}{\sqrt{n}} \leq \mu \leq \bar{X} + \frac{z_{0.025}}{\sqrt{n}}\right) = 95\%.$$

We get a **random** interval $\left[\bar{X} - \frac{z_{0.025}}{\sqrt{n}}, \bar{X} + \frac{z_{0.025}}{\sqrt{n}}\right]$ which can cover $\mu$ with 95% probability!



*: The "realization" of the random interval that does NOT cover $\mu$

# Revisit the derivation of CI of $\mu$ in $N(\mu, 1)$

**Steps to construct a CI**: Suppose we have samples $X_1, \ldots, X_n$, and $\theta$ is the unknown parameter.

- **Step 1:** Find a r.v. $V$ which is a function of both $X_1, \ldots, X_n$ and $\theta$
- **Step 2:** Verify that the distribution of $V$ does NOT depend on $\theta$ or any other unknown parameters
- **Step 3:** Derive the equation from the fact that
  $\mathbb{P}(v_{1-\alpha/2} \leq V(X_1, \ldots, X_n, \theta) \leq v_{\alpha/2}) = 1 - \alpha$, where $\mathbb{P}(V \geq v_\beta) = \beta$ for any $\beta \in [0, 1]$.

Suppose $X_1, \ldots, X_n \overset{i.i.d.}{\sim} N(\mu, 1)$. $\bar{X} = \frac{1}{n} \sum_{i=1}^{n} X_i \sim N(\mu, 1/n)$. By standardization, $\frac{\bar{X} - \mu}{\sqrt{1/n}} \sim N(0, 1)$. Thus,

$$\mathbb{P}\left( z_{0.975} \leq \frac{\bar{X} - \mu}{\sqrt{1/n}} \leq z_{0.025} \right) = 95\%,$$
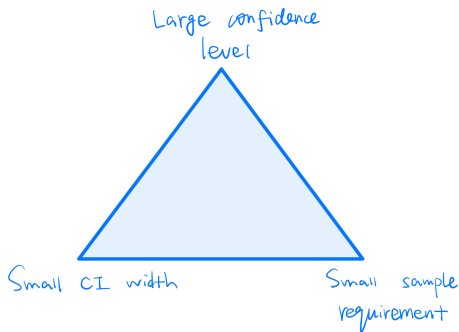
Then

$$\mathbb{P}\left( \bar{X} - \frac{z_{0.025}}{\sqrt{n}} \leq \mu \leq \bar{X} - \frac{z_{0.975}}{\sqrt{n}} \right) = 95\%.$$

# Confidence level, precision (width), and sample size

The ideal scenario:
- **Large** confidence level
- **Small** CI width
- **Small** sample size requirement

But this is an **impossible trinity**!



Large confidence level

Small CI width

Small sample requirement

$$\mathbb{P}\left( \bar{X} - \frac{z_{\alpha/2}}{\sqrt{n}} \leq \mu \leq \bar{X} + \frac{z_{\alpha/2}}{\sqrt{n}} \right) = 1 - \alpha.$$

Confidence level $= 100(1 - \alpha)\%$, CI width $= \frac{2z_{\alpha/2}}{\sqrt{n}}$, sample size $n$.
- Fixed confidence level: CI width $\searrow$, $n$ needs to be $\nearrow$
- Fixed confidence level: $n \searrow$, CI width $\nearrow$
- Fixed sample size $n$: CI width $\nearrow$, corresponding confidence level $\nearrow$
- Fixed sample size CI width: $n \nearrow$, corresponding confidence level $\nearrow$

# $\chi^2$-distribution and t-distribution

$\chi^2$-**distribution:** Suppose $Z_1, \ldots, Z_p \overset{i.i.d.}{\sim} N(0,1)$. Then we say variable $V = \sum_{i=1}^{p} Z_i^2$ follows the $\chi^2$-**distribution** with **degree of freedom** $p$, denoted as $V \sim \chi_p^2$.

t-**distribution:** Suppose $Z \sim N(0,1)$, $Q \sim \chi_p^2$, and $Z$ is independent with $Q$. Then we say the variable $\frac{Z}{\sqrt{Q/p}}$ follows **t-distribution with degree of freedom** $p$, denoted as $\frac{Z}{\sqrt{Q/p}} \sim t_p$.

**Applications:** Suppose $X_1, \ldots, X_n \overset{i.i.d.}{\sim} N(\mu, \sigma^2)$. The sample variance $S^2 = \frac{1}{n-1} \sum_{i=1}^{n} (X_i - \bar{X})^2$. Then:

○ $\frac{(n-1)s^2}{\sigma^2} \sim \chi_{n-1}^2$.

○ $\frac{\sqrt{n}(\bar{X} - \mu)}{S} \sim t_{n-1}$.

# Today's goal

○ Know the CIs of some commonly used models and understand how to derive them.

    ▷ Normal distribution $N(\mu, \sigma^2)$

        ◇ $\sigma^2$ known, derive CI of $\mu$

        ◇ $\sigma^2$ unknown, derive CI of $\mu$

        ◇ Derive CI of $\sigma^2$ (or $\sigma$)

    ▷ Bernoulli distribution Bernoulli($p$): derive CI of $p$

    ▷ Exponential distribution Exp($\lambda$): derive CI of $\lambda$ (you will do it in HW4)

○ Understand the relation (impossible trinity) between confidence level, CI width and minimum sample size requirement, and can interpret it in specific models.

# CI in Normal Distribution $N(\mu, \sigma^2)$

# $\sigma^2$ known, derive CI of $\mu$

**Steps to construct a CI**:

- **Step 1:** Find a r.v. $V$ which is a function of both $X_1, \ldots, X_n$ and $\mu$
- **Step 2:** Verify that the distribution of $V$ does NOT depend on $\mu$ or any other unknown parameters
- **Step 3:** Derive the equation from the fact that $\mathbb{P}(v_{1-\alpha/2} \leq V(X_1, \ldots, X_n, \mu) \leq v_{\alpha/2}) = 1 - \alpha$, where $\mathbb{P}(V \geq v_\beta) = \beta$ for any $\beta \in [0, 1]$.

Suppose $X_1, \ldots, X_n \overset{i.i.d.}{\sim} N(\mu, \sigma^2)$. $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \sim N(\mu, \sigma^2/n)$. By standardization, $\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \sim N(0, 1)$. Thus,

$$\mathbb{P}\left(-z_{\alpha/2} \leq \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \leq z_{\alpha/2}\right) = 1 - \alpha,$$

Then

$$\mathbb{P}\left(\bar{X} - \frac{\sigma}{\sqrt{n}} z_{\alpha/2} \leq \mu \leq \bar{X} + \frac{\sigma}{\sqrt{n}} z_{\alpha/2}\right) = 1 - \alpha.$$

Thus a $100(1-\alpha)\%$ CI of $\mu$: $\left[\bar{X} - \frac{\sigma}{\sqrt{n}} z_{\alpha/2}, \bar{X} + \frac{\sigma}{\sqrt{n}} z_{\alpha/2}\right]$.

# Example: $\sigma^2$ known, derive CI of $\mu$

**Example:** Extensive monitoring of a computer time-sharing system has suggested that response time to a particular editing command is normally distributed with standard deviation 25 millisec. We wish to estimate the true average response time $\mu$. What sample size is necessary to ensure that the resulting 95% CI has a width of (at most) 10?

Recall that the 95% CI is $\left[\bar{X} - \frac{\sigma}{\sqrt{n}} z_{0.025}, \bar{X} + \frac{\sigma}{\sqrt{n}} z_{0.025}\right] =$
$\left[\bar{X} - \frac{25 \times 1.96}{\sqrt{n}}, \bar{X} + \frac{25 \times 1.96}{\sqrt{n}}\right]$. The width $= 2 \times \frac{25 \times 1.96}{\sqrt{n}}$. Thus we need

$$2 \times \frac{25 \times 1.96}{\sqrt{n}} \leq 10 \Rightarrow n \geq \left(\frac{2 \times 25 \times 1.96}{10}\right)^2 = 96.04.$$

Therefore, we need sample size $n \geq 97$.

# Derive CI of $\sigma^2$ (or $\sigma$)

**Steps to construct a CI**:

○ **Step 1:** Find a r.v. $V$ which is a function of both $X_1, \ldots, X_n$ and $\sigma^2$

○ **Step 2:** Verify that the distribution of $V$ does NOT depend on $\sigma^2$ or any other unknown parameters

○ **Step 3:** Derive the equation from the fact that
$\mathbb{P}(v_{1-\alpha/2} \leq V(X_1, \ldots, X_n, \sigma^2) \leq v_{\alpha/2}) = 1 - \alpha$, where $\mathbb{P}(V \geq v_\beta) = \beta$ for any $\beta \in [0,1]$.

$X_1, \ldots, X_n \overset{i.i.d.}{\sim} N(\mu, \sigma^2)$. Sample variance $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$.

Then $\frac{(n-1)S^2}{\sigma^2} \sim \chi^2_{n-1}$, leading to

$$\mathbb{P}\left(\chi^2_{n-1,1-\alpha/2} \leq \frac{(n-1)S^2}{\sigma^2} \leq \chi^2_{n-1,\alpha/2}\right) = 1 - \alpha,$$

which implies

$$\mathbb{P}\left(\frac{(n-1)S^2}{\chi^2_{n-1,\alpha/2}} \leq \sigma^2 \leq \frac{(n-1)S^2}{\chi^2_{n-1,1-\alpha/2}}\right) = 1 - \alpha.$$

So a $100(1-\alpha)\%$ CI of $\sigma^2$: $\left[\frac{(n-1)S^2}{\chi^2_{n-1,\alpha/2}}, \frac{(n-1)S^2}{\chi^2_{n-1,1-\alpha/2}}\right]$.

# Derive CI of $\sigma^2$ (or $\sigma$)

$X_1, \ldots, X_n \overset{i.i.d.}{\sim} N(\mu, \sigma^2)$. Sample variance $S^2 = \frac{1}{n-1} \sum_{i=1}^{n} (X_i - \bar{X})^2$.
A $100(1-\alpha)\%$ CI of $\sigma^2$: $\left[ \frac{(n-1)S^2}{\chi^2_{n-1,\alpha/2}}, \frac{(n-1)S^2}{\chi^2_{n-1,1-\alpha/2}} \right]$.

$$\mathbb{P}\left( \frac{(n-1)S^2}{\chi^2_{n-1,\alpha/2}} \leq \sigma^2 \leq \frac{(n-1)S^2}{\chi^2_{n-1,1-\alpha/2}} \right) = 1 - \alpha.$$

This implies

$$\mathbb{P}\left( \sqrt{\frac{(n-1)S^2}{\chi^2_{n-1,\alpha/2}}} \leq \sigma \leq \sqrt{\frac{(n-1)S^2}{\chi^2_{n-1,1-\alpha/2}}} \right) = 1 - \alpha.$$

A $100(1-\alpha)\%$ CI of $\sigma$: $\left[ \sqrt{\frac{(n-1)S^2}{\chi^2_{n-1,\alpha/2}}}, \sqrt{\frac{(n-1)S^2}{\chi^2_{n-1,1-\alpha/2}}} \right]$.

This idea can be generalized into the following result: For a $100(1-\alpha)\%$ CI of $\theta$ (denoted as $[l, u]$), the $100(1-\alpha)\%$ CI of $g(\theta)$ equals

○ $[g(l), g(u)]$, if $g$ is an increasing function

○ $[g(u), g(l)]$, if $g$ is a decreasing function

# Example: Derive CI of $\sigma^2$ (or $\sigma$)

The accompanying data on breakdown voltage of electrically stressed circuits was read from a normal probability plot. The straightness of the plot gave strong support to the assumption that breakdown voltage is approximately normally distributed.

1470, 1510, 1690, 1740, 1900, 2000, 2030, 2100, 2190, 2200, 2290, 2380, 2390, 2480, 2500, 2580, 2700

Let $\sigma^2$ denote the variance of the breakdown voltage distribution. Derive a 90% CI for $\sigma$.

---

$n = 17 \Rightarrow$ df $= n - 1 = 16$. The sample variance $s^2 = 137324.3$. $\chi^2_{16,0.95} = 7.962$ and $\chi^2_{16,0.05} = 26.296$.

Recall that a 90% CI of $\sigma^2$ is

$\left[ \frac{(n-1)s^2}{\chi^2_{n-1,0.05}}, \frac{(n-1)s^2}{\chi^2_{n-1,0.95}} \right] = \left[ \frac{16 \times 137324.3}{26.296}, \frac{16 \times 137324.3}{7.962} \right] = [83556.01, 275959.4]$.

Thus a 90% CI of $\sigma$ is $\left[ \sqrt{83556.01}, \sqrt{275959.4} \right] = [289.06, 525.32]$.

# $\sigma^2$ unknown, derive CI of $\mu$

**Steps to construct a CI**:

◦ **Step 1:** Find a r.v. $V$ which is a function of both $X_1, \ldots, X_n$ and $\mu$

◦ **Step 2:** Verify that the distribution of $V$ does NOT depend on $\mu$ or any other unknown parameters

◦ **Step 3:** Derive the equation from the fact that
$\mathbb{P}(v_{1-\alpha/2} \leq V(X_1, \ldots, X_n, \mu) \leq v_{\alpha/2}) = 1 - \alpha$, where $\mathbb{P}(V \geq v_\beta) = \beta$ for any $\beta \in [0, 1]$.

Suppose $X_1, \ldots, X_n \overset{i.i.d.}{\sim} N(\mu, \sigma^2)$. $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \sim N(\mu, \sigma^2/n)$. By standardization, $\frac{\bar{X}-\mu}{\sigma/\sqrt{n}} \sim N(0, 1)$. Thus,

$$\mathbb{P}\left( \bar{X} - \frac{\sigma}{\sqrt{n}} z_{\alpha/2} \leq \mu \leq \bar{X} + \frac{\sigma}{\sqrt{n}} z_{\alpha/2} \right) = 1 - \alpha.$$

Can we still claim a $100(1-\alpha)\%$ CI of $\mu$ is $\left[ \bar{X} - \frac{\sigma}{\sqrt{n}} z_{\alpha/2}, \bar{X} + \frac{\sigma}{\sqrt{n}} z_{\alpha/2} \right]$?
**NO!!!** Because $\sigma$ is **unknown**!

# $\sigma^2$ unknown, derive CI of $\mu$

Suppose $X_1, \ldots, X_n \overset{i.i.d.}{\sim} N(\mu, \sigma^2)$. $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \sim N(\mu, \sigma^2/n)$.
Sample variance $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$.
From the last lecture, we know that

○ $\frac{\bar{X}-\mu}{\sigma/\sqrt{n}} \sim N(0,1)$;

○ $\frac{(n-1)S^2}{\sigma^2} \sim \chi^2_{n-1}$;

○ $\bar{X}$ and $S^2$ are independent

Thus, by definition of t-distribution, $\frac{\sqrt{n}(\bar{X}-\mu)}{S} \sim t_{n-1}$, implying that

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

Then

$$\mathbb{P}\left( \bar{X} - \frac{S}{\sqrt{n}} t_{n-1,\alpha/2} \leq \mu \leq \bar{X} + \frac{S}{\sqrt{n}} t_{n-1,\alpha/2} \right) = 1 - \alpha.$$

Thus a $100(1-\alpha)\%$ CI of $\mu$: $\left[ \bar{X} - \frac{S}{\sqrt{n}} t_{n-1,\alpha/2}, \bar{X} + \frac{S}{\sqrt{n}} t_{n-1,\alpha/2} \right]$.

# Compare two cases: $\sigma^2$ known/unknown

$\sigma^2$ **known:** A $100(1-\alpha)\%$ CI of $\mu$: $\left[\bar{X} - \frac{\sigma}{\sqrt{n}}z_{\alpha/2}, \bar{X} + \frac{\sigma}{\sqrt{n}}z_{\alpha/2}\right]$.

$\sigma^2$ **unknown:** A $100(1-\alpha)\%$ CI of $\mu$: $\left[\bar{X} - \frac{S}{\sqrt{n}}t_{n-1,\alpha/2}, \bar{X} + \frac{S}{\sqrt{n}}t_{n-1,\alpha/2}\right]$.

The only two differences from $\sigma^2$ known to $\sigma^2$ unknown:

∘ Replacing true SD $\sigma$ by sample SD $S$;

∘ Replacing z-values by t-values.

Connection:

∘ By Weak Law of Large Number, $S \approx \sigma$ when $n$ is large

∘ When the degree of freedom is large (here it means $n$ is large), t-distribution becomes very similar to standard normal distribution
$\Rightarrow$ when $n \geq 30$, can replace $t_{n-1,\alpha/2}$ by $z_{\alpha/2}$

# Example: $\sigma^2$ unknown, derive CI of $\mu$

An object is weighed 9 times, with an average weight 1.03 kg and SD 0.10 kg. Calculate the $95\%$ CI for the unknown weight.

---

$\sigma^2$ **unknown** and $n = 9 < 30$, therefore we will use t-distribution and the corresponding quantiles.

We know that $\bar{x} = 1.03$, $s = 0.10$.

A $100(1-\alpha)\%$ CI of $\mu$: $\left[ \bar{x} - \frac{s}{\sqrt{9}} t_{8,0.025}, \bar{x} + \frac{s}{\sqrt{n}} t_{8,0.025} \right] =$

$\left[ 1.03 - \frac{0.1}{3} \times 2.306, 1.03 + \frac{0.1}{3} \times 2.306 \right] = [0.953, 1.107]$.

# CI of $p$ in Bernoulli$(p)$

# CI of success probability $p$

**Example:** If the sample includes 100 employees, find a 95% confidence interval for **the proportion** of employees who don't like their jobs in the sample.

Therefore sometimes we also call it the CI for a **population proportion** $p$.

# Derive the CI of success probability $p$

**Steps to construct a CI**: Suppose we have samples $X_1, \ldots, X_n$, and $\theta$ is the unknown parameter.

- **Step 1:** Find a r.v. $V$ which is a function of both $X_1, \ldots, X_n$ and $\theta$
- **Step 2:** Verify that the distribution of $V$ does NOT depend on $\theta$ or any other unknown parameters
- **Step 3:** Derive the equation from the fact that
  $\mathbb{P}(v_{1-\alpha/2} \leq V(X_1, \ldots, X_n, \theta) \leq v_{\alpha/2}) = 1 - \alpha$, where $\mathbb{P}(V \geq v_\beta) = \beta$ for any $\beta \in [0, 1]$.

Suppose $X_1, \ldots, X_n \overset{i.i.d.}{\sim}$ Bernoulli$(p)$. By central limit theorem $\hat{p} = \bar{X} = \frac{1}{n}\sum_{i=1}^{n} X_i \overset{d}{\approx} N(p, p(1-p)/n)$. By standardization, $\frac{\hat{p}-p}{\sqrt{p(1-p)/n}} \overset{d}{\approx} N(0, 1)$. Thus,

$$\mathbb{P}\left( z_{0.975} \leq \frac{\hat{p} - p}{\sqrt{p(1-p)/n}} \leq z_{0.025} \right) \approx 95\%,$$

Then

$$\mathbb{P}\left( z_{0.975} \leq \frac{\hat{p} - p}{\sqrt{\hat{p}(1-\hat{p})/n}} \leq z_{0.025} \right) \approx 95\%,$$

# Derive the CI of success probability $p$

Suppose $X_1, \ldots, X_n \overset{i.i.d.}{\sim}$ Bernoulli($p$). By central limit theorem $\hat{p} = \bar{X} = \frac{1}{n} \sum_{i=1}^{n} X_i \overset{d}{\approx} N(p, p(1-p)/n)$. By standardization, $\frac{\hat{p}-p}{\sqrt{p(1-p)/n}} \sim N(0,1)$. Thus,

$$\mathbb{P}\left( z_{0.975} \leq \frac{\hat{p}-p}{\sqrt{p(1-p)/n}} \leq z_{0.025} \right) \approx 95\%.$$

Then

$$\mathbb{P}\left( z_{0.975} \leq \frac{\hat{p}-p}{\sqrt{\hat{p}(1-\hat{p})/n}} \leq z_{0.025} \right) \approx 95\%.$$

$$\mathbb{P}\left( \hat{p} - z_{0.025}\sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \leq p \leq \hat{p} + z_{0.025}\sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \right) \approx 95\%.$$

An **approximate** $100(1-\alpha)\%$ CI of $p$:

$$\left[ \hat{p} - z_{0.025}\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}, \hat{p} + z_{0.025}\sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \right]$$

# Example: CI of success probability $p$

Construct a 95% CI for the true proportion of college students who sleep fewer than 6 hours per night, if the sample proportion in a sample of 500 students is 0.3.

---

$\hat{p} = 0.3$, $n = 500$. Plug them into our formula:

An approximate $100(1-\alpha)\%$ CI of $p$ is

$$\left[\hat{p} - z_{0.025}\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}, \hat{p} + z_{0.025}\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}\right] =$$

$$\left[0.3 - 1.96\sqrt{\frac{0.3 \times 0.7}{500}}, 0.3 + 1.96\sqrt{\frac{0.3 \times 0.7}{500}}\right] = [0.2598, 0.3402]$$

# Example: sample size requirement

A college dean wishes to survey the undergraduate population to find out what proportion $p$ of the students would prefer to eliminate all 8:40am course offerings. What sample size is needed if the 95% CI for $p$ is to have a width of at most 0.06 **irrespective of** $\hat{p}$?

---

Actually before we do the survey, we don't know $\hat{p} \Rightarrow$ We have to figure out the minimum sample size that works for **every** possible values of $\hat{p} \in [0, 1]$.

Recall: A $100(1 - \alpha)\%$ CI of $p$: $\left[ \hat{p} - z_{0.025}\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}, \hat{p} + z_{0.025}\sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \right]$

We want CI width $= 2z_{0.025}\sqrt{\frac{\hat{p}(1-\hat{p})}{n}} = 2 \times 1.96\sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \leq 0.06$ holds for all $\hat{p} \in [0, 1] \Rightarrow n \geq (\frac{2 \times 1.96}{0.06})^2 \hat{p}(1-\hat{p})$ for all $\hat{p} \in [0, 1]$.
By Cauchy-Schwarz inequality (HW0, or directly maximization of quadratic function $\hat{p}(1 - \hat{p})$): $\hat{p}(1 - \hat{p}) \leq (\frac{\hat{p}+1-\hat{p}}{2})^2 = 0.25$. Thus if suffices to have $n \geq (\frac{2 \times 1.96}{0.06})^2 \times 0.25 = 1067.1 \Rightarrow n_{\min} = 1068$.

# Example: Wordle



| Date | No. | Word | Yongxin | Ye |
|---|---|---|---|---|
| 1/18/22 | 213 | proxy | 6 | 4 |
| 1/19/22 | 214 | point | 4 | 4 |
| 1/20/22 | 215 | robot | 4 | 5 |
| 1/21/22 | 216 | prick | 6 | 4 |
| 1/22/22 | 217 | wince | 5 | 4 |

# Example: Wordle

Suppose the number of tries of Xin and Ye follow some normal distribution $N(\mu_1, \sigma^2)$ and $N(\mu_2, \sigma^2)$, which shares the same unknown variance $\sigma^2$. Construct a 95% CI of $\mu_1 - \mu_2$.

Here're the number of guesses from Jan 18, 2022 to Feb 18, 2022:

○ Xin ($X_i$): 6 4 4 6 5 4 6 5 5 4 4 4 4 3 3 3 3 4 5 3 6 5 5 4 5 5 4 5 4 5 3 5

○ Ye ($Y_i$):   4 4 5 4 4 4 5 5 6 5 6 2 5 6 3 3 6 4 6 4 6 5 4 2 4 3 5 4 3 3 3 6

○ Difference ($D_i = X_i - Y_i$): 2 0 -1 2 1 0 1 0 -1 -1 -2 2 -1 -4 0 0 -3 0 -1 -1 0 0 1 2 1 2 -1 1 1 2 0 -1

---

$D_i = X_i - Y_i$'s are independent from each other, $D_i \sim N(\mu_1 - \mu_2, 2\sigma^2)$

$\bar{D} = \frac{1}{n}\sum_{i=1}^{n} D_i \sim N(\mu_1 - \mu_2, 2\sigma^2/n), \quad n = 32$

Sample mean $\bar{d} = 0.0625$.

Sample variance of the difference equals $s^2 = \frac{1}{n-1}\sum_{i=1}^{n}(D_i - \bar{D})^2$

$= 1.9315 \quad \Rightarrow \quad s = 1.3898$

A $100(1-\alpha)\%$ CI of $\mu_1 - \mu_2$: $\left[\bar{D} - \frac{S}{\sqrt{n}}t_{31,0.025}, \bar{D} + \frac{S}{\sqrt{n}}t_{31,0.025}\right]$

$= [0.0625 - \frac{1.3898}{\sqrt{32}} \times 2.04, 0.0625 + \frac{1.3898}{\sqrt{32}} \times 2.04] = [-0.4387, 0.5637]$.

# Example: Wordle

Suppose the number of tries of Xin and Ye follow some normal distribution $N(\mu_1, \sigma^2)$ and $N(\mu_2, \sigma^2)$, which shares the same unknown variance $\sigma^2$. Construct a 95% CI of $\mu_1 - \mu_2$.

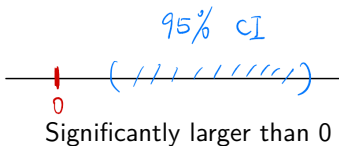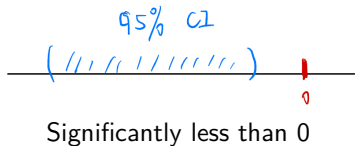Here're the number of guesses from Jan 18, 2022 to Feb 18, 2022:

- Xin ($X_i$): 6 4 4 6 5 4 6 5 5 4 4 4 4 3 3 3 3 4 5 3 6 5 5 4 5 5 4 5 4 5 3 5
- Ye ($Y_i$):   4 4 5 4 4 4 4 5 5 6 5 6 2 5 6 3 3 6 4 6 4 6 5 4 2 4 3 5 4 3 3 3 6
- Difference ($D_i = X_i - Y_i$): 2 0 -1 2 1 0 1 0 -1 -1 -2 2 -1 -4 0 0 -3 0 -1 -1 0 0 1 2 1 2 -1 1 1 2 0 -1

A $100(1-\alpha)\%$ CI of $\mu_1 - \mu_2$: $\left[\bar{D} - \frac{S}{\sqrt{n}}t_{31,0.025}, \bar{D} + \frac{S}{\sqrt{n}}t_{31,0.025}\right]$
$= [0.0625 - \frac{1.3898}{\sqrt{32}} \times 2.04, 0.0625 + \frac{1.3898}{\sqrt{32}} \times 2.04] = [-0.4387, 0.5637]$.

Do you think any one does significantly better than the other?

NO! Because the 95% CI covers 0!



Significantly less than 0        Significantly larger than 0

**Many thanks to**

- Joyce Robbins
- Yang Feng
- Chengliang Tang
- Owen Ward
- Wenda Zhou
- Yongxin Shang
- And all my teachers in the past 25 years