

CSCI 3022

intro to data science with probability & statistics

Lecture 15
March 7, 2018

Introduction to Statistical Inference & Confidence Intervals

NOTE: It was L'Hospital who happily stole Johann Bernoulli's work & published it as his own. The next time you take an indeterminate limit, remember that & call it "Bernoulli's Rule" instead!

TONY
WUZ HERE

Last time on CSCI 3022

- **Proposition:** If X is a normally distributed random variable with mean μ and standard deviation σ , then Z is a standard normal distribution if

$$Z = \frac{X - \mu}{\sigma} \quad \text{or} \quad X = \sigma Z + \mu$$

- **The Central Limit Theorem:** Let X_1, X_2, \dots, X_n be i.i.d. draws from some distribution. Then as n becomes large

non-normal
or
normal

$$\bar{X} \sim N\left(\mu, \frac{\sigma^2}{n}\right) \quad \text{2}$$

- A $100(1 - \alpha)\%$ confidence interval for the mean μ when the value of σ is known is given by:

↳
$$\left[\bar{X} - z_{\alpha/2} \frac{\sigma}{\sqrt{n}}, \bar{X} + z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \right]$$

Statistical Inference

- **Goal:** Want to extract properties of an underlying population by analyzing sampled data
- **Last time we saw:**
 - How to determine a confidence interval for the population mean
 - How to determine a confidence interval for the population proportion
- **This time we'll see:**
 - How to put a confidence interval on the difference between means of two populations
 - How to put a confidence interval on the difference between proportions of two populations
 - How we can get a good numerical estimate of a CI using something called the Bootstrap

Difference between population means

- How do two sub-populations compare? In particular, are their means the same?
- **Classic Motivating Examples:**
 - Is a drug's effectiveness the same in children and adults? ✓
 - Does cigarette brand A contain more nicotine than cigarette brand B?
 - Does a class perform better when Professor C ^{teaches} it or Professor D ^{teaches} it?
 - Does email ad E generate more customers than email ad F?

Difference between population means

- How do two sub-populations compare? In particular, are their means the same?
- **Solution Process:**
 - Collect samples from both sub-populations, and perform inference on both samples to make conclusions about $\mu_1 - \mu_2$

Difference between population means

- How do two sub-populations compare? In particular, are their means the same?
- **Solution Process:**
 - Collect samples from both sub-populations, and perform inference on both samples to make conclusions about $\mu_1 - \mu_2$

- **Basic Assumptions:**

- OLD NEWS*
- (X_1, X_2, \dots, X_m) is a random sample from a distribution with mean μ_1 and sd σ_1
m X's
 - (Y_1, Y_2, \dots, Y_n) is a random sample from a distribution with mean μ_2 and sd σ_2
n Y's
 - The X and Y samples are independent of each other.
- indep.*
- \bar{x}*
 s_x

Difference between population means

- The natural estimator of $\mu_1 - \mu_2$ is the difference of the sample means $\bar{x} - \bar{y}$
- Is $\bar{x} - \bar{y}$ a good estimator for $\mu_1 - \mu_2$?

$\bar{X} - \bar{Y}$ s.v.

- The expected value of $\bar{X} - \bar{Y}$ is given by

$$\begin{aligned} E[\bar{X} - \bar{Y}] &= E[\bar{X}] - E[\bar{Y}] \\ &= \underbrace{\mu_1} - \underbrace{\mu_2} \quad \checkmark \checkmark \end{aligned}$$

- The standard deviation of $\bar{X} - \bar{Y}$ is given by

$$\begin{aligned} SD[\bar{X} - \bar{Y}] &= \sqrt{\text{Var}[\bar{X} - \bar{Y}]} = \sqrt{\text{Var}[\bar{X}] + \text{Var}[\bar{Y}]} = \sqrt{\frac{\sigma_1^2}{n} + \frac{\sigma_2^2}{n}} \quad \checkmark \end{aligned}$$

Handwritten notes: Arrows indicate the derivation from $\text{Var}[\bar{X} - \bar{Y}]$ to $\text{Var}[\bar{X}] + \text{Var}[\bar{Y}]$ using the property $(-1)^2 = 1$, and then to the final formula.

Normal populations with known SDs

- If both populations are normal, then both \bar{X} and \bar{Y} are normally distributed.
- Independence of the two samples implies that the sample means are independent.
- Therefore, the difference between the means is normally distributed, for any sample sizes, with:

$$\underbrace{\bar{X} - \bar{Y}}_{\text{estimator}} \sim N\left(\underbrace{\mu_1 - \mu_2}_{\substack{\text{expected} \\ \text{value} \\ \text{of} \\ \text{est}_1}}, \underbrace{\frac{\sigma_1^2}{m} + \frac{\sigma_2^2}{n}}_{\text{variance}}\right)$$

Confidence Interval for the difference

- Standardizing $\bar{X} - \bar{Y}$ gives a standard normal random variable

$$Z = \frac{(\bar{x} - \bar{y}) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{m} + \frac{\sigma_2^2}{n}}}$$

messy?

But oh so nice!

$$\sim \underline{N(0, 1)}$$

- And so, we can compute a 100(1 - α)% confidence interval for $\mu_1 - \mu_2$

$$\bar{x} - \bar{y} \pm z_{\alpha/2} \sqrt{\frac{\sigma_1^2}{m} + \frac{\sigma_2^2}{n}}$$

central est. \pm z-score * SD

Large sample CIs for the difference

- **Not surprisingly**, if both m and n are large, then our friend, the CLT, kicks in, and our confidence interval for the difference of means is valid, even when the populations are *not* normally distributed!
- **Furthermore**, if m and n are large, and we don't know the standard deviations, we can replace them with the sample standard deviations:

$$\sigma_1^2 \mapsto s_1^2 = \frac{1}{m-1} \sum_i (x_i - \bar{x})^2$$

$$\sigma_2^2 \mapsto s_2^2 = \frac{1}{n-1} \sum_j (y_j - \bar{y})^2$$

Confidence Interval for the Difference

$$[-0.508, 0.068]$$

- Example:** Suppose you run two different email Ad campaigns over many days and record the amount of traffic driven to your website on days that each ad is sent. In particular, suppose that Ad 1 is sent on 50 different days and generates an average of 2 million page views per day with an sd of 1 million views, and Ad 2 is sent on 40 different days and generates an average of 2.25 million page views per day with an sd of a half million views. Find a 95% confidence interval for the difference in average page views per day (in units of millions of views).

$$\begin{aligned}\bar{x} &= 2 \\ n &= 50 \\ s_1 &= 1\end{aligned}$$

$$\begin{aligned}z_{\alpha/2} &= z_{0.05/2} \\ &= 1.96\end{aligned}$$

$$\begin{aligned}\bar{y} &= 2.25 \\ n &= 40 \\ s_2 &= 0.5\end{aligned}$$

$$95\% \text{ CI} = (\bar{x} - \bar{y}) \pm z_{\alpha/2} \sqrt{\frac{s_1^2}{n} + \frac{s_2^2}{n}}$$

$$= (2 - 2.25) \pm 1.96 \cdot \sqrt{\frac{1^2}{50} + \frac{0.5^2}{40}} = [-0.508, 0.068]$$

Confidence Interval for the Difference

- **Example:** Suppose you run two different email Ad campaigns over many days and record the amount of traffic driven to your website on days that each ad is sent. In particular, suppose that Ad 1 is sent on 50 different days and generates an average of 2 million page views per day with an sd of 1 million views, and Ad 2 is sent on 40 different days and generates an average of 2.25 million page views per day with an sd of a half million views. Find a 95% confidence interval for the difference in average page views per day (in units of millions of views).

Confidence Interval for the Difference

HERE



$$CI \text{ width: } 2 \cdot z_{\alpha/2} \cdot SD$$

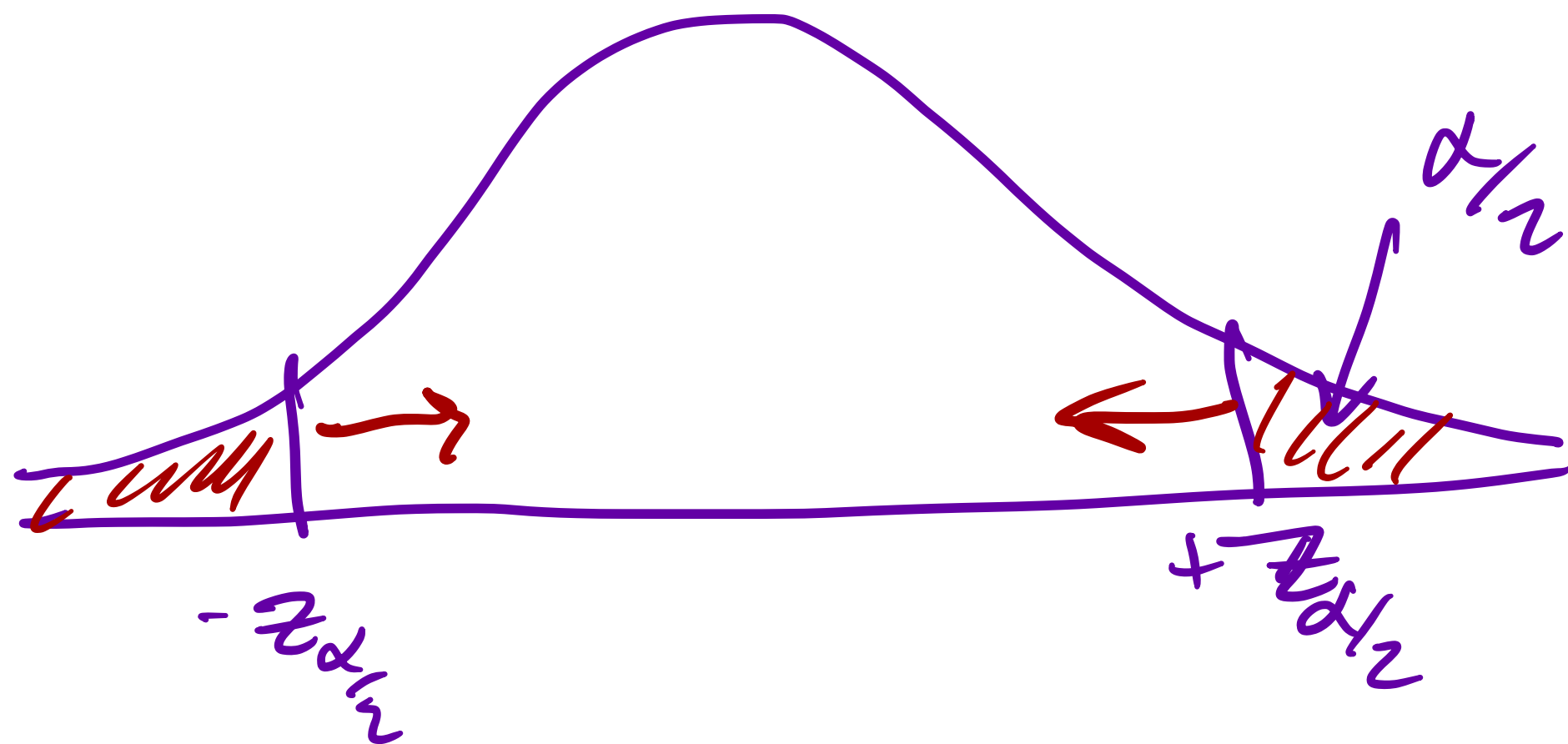
- **Looking forward to interpretation:** What does our confidence interval tell us about the effectiveness of the two advertisements?

$$[-0.5, +0.1] \text{ (ish)}$$

$$\bar{X} < \bar{Y} \rightarrow \bar{Y} \text{ is better?}$$

contains 0 \rightarrow

so no statistically significant difference
at $\alpha = 0.05$ confidence level



What happens if we increase α ?

$$\alpha \uparrow \Rightarrow z_{\alpha/2} \downarrow \Rightarrow \text{CI width} \downarrow \Rightarrow 0 \text{ gets kicked out}$$