

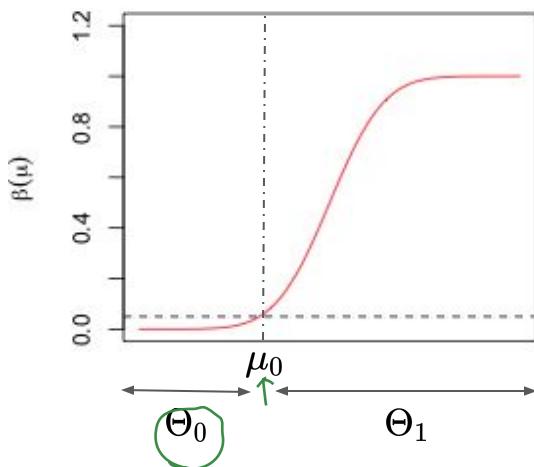
Generalized LRT

Examples

07/13/2021

In the previous lecture,

$$\underline{H_0 : \mu \leq \mu_0} \quad \text{versus} \quad H_1 : \mu > \mu_0.$$



- LRT is uniformly most powerful:
 - Neyman-Pearson Lemma:
 $H_0 : \theta = \theta_0 \leftrightarrow H_1 : \theta = \theta_1 \ (\theta_1 > \theta_0).$
 - Proof in class:
 $H_0 : \theta = \theta_0 \leftrightarrow H_1 : \theta > \theta_0.$ ↘
 - Karlin-Rubin Theorem under monotonicity:
 $H_0 : \theta \leq \theta_0 \leftrightarrow H_1 : \theta > \theta_0.$ ↘
- LRT under $N(\mu, \sigma^2)$:
 - σ^2 known or unknown;
 - Duality between CIs and hypothesis tests
- Generalized LRT:
 - When the exact sampling distribution of $\lambda(\mathbf{X}_n)$ is hard to obtain;
 - Under H_0 , Wilk's theorem guarantees

$$-2 \log \lambda(\mathbf{X}_n) \xrightarrow{d} \chi_{\nu}^2 \text{ as } n \rightarrow \infty.$$

~~~~~

$$P(X=x | \lambda) = \frac{e^{-\lambda} \lambda^x}{x_i!}$$

$$\hat{\lambda}_{MLE} = \bar{x}_n$$

## Generalized LRT example

**Example 6.** Let  $X_1, \dots, X_{25}$  be i.i.d  $Poisson(\lambda)$ . Consider

$$H_0 : \lambda = 5 \leftrightarrow H_1 : \lambda \neq 5.$$

$$\{\lambda(\bar{x}_n) \leq c\}$$

Solution:  $L(\lambda | \bar{x}_n) = \frac{e^{-n\lambda} \lambda^{\sum x_i}}{\prod_{i=1}^{25} x_i!} e^{\sum x_i \log 5} \quad H_0 = \{\lambda = 5\}$

$$\lambda(\bar{x}_n) = \frac{\sup_{H_0} L(\lambda | \bar{x}_n)}{\sup_{H_1} L(\lambda | \bar{x}_n)} = \frac{e^{-n\bar{x}_n} 5^{\sum x_i}}{e^{-n\bar{x}_n} \bar{x}_n^{\sum x_i} / \prod_{i=1}^{25} x_i!} = e^{n(\bar{x}_n - 5)} e^{(\sum x_i)(\log 5 - \log \bar{x}_n)}$$

$\uparrow$  unrestricted

$$e^{(\sum x_i) \log \bar{x}_n}$$

$$= e^{n(\bar{x}_n - 5) + n \bar{x}_n (\log 5 - \log \bar{x}_n)}$$

only depends on  $\bar{x}_n$ .

$$-2 \log \lambda(\bar{x}_n) = -2 [n(\bar{x}_n - 5) + n \bar{x}_n (\log 5 - \log \bar{x}_n)] \xrightarrow{d} \chi_1^2.$$

# Generalized LRT example

**Example 6.** Let  $X_1, \dots, X_{25}$  be i.i.d  $\text{Poisson}(\lambda)$ . Consider

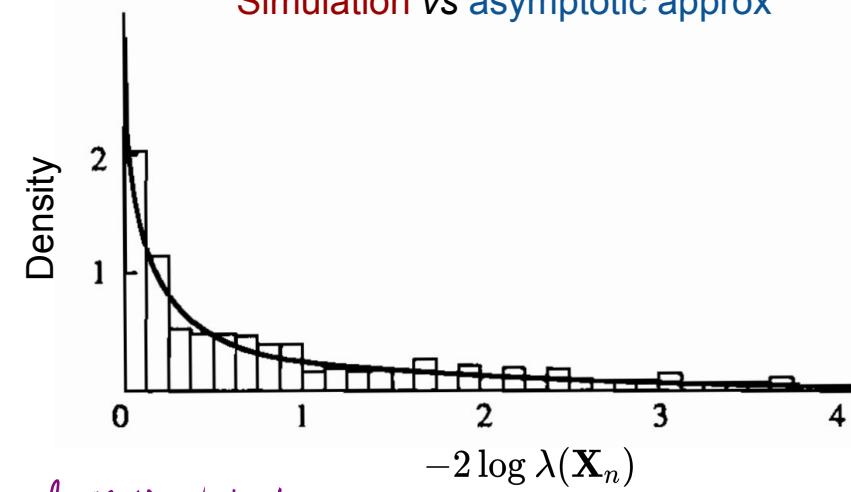
$$H_0 : \lambda = 5 \leftrightarrow H_1 : \lambda \neq 5.$$

$$-2 \log \lambda(\mathbf{X}_n) = -2n[(\bar{X}_n - 5) + \bar{X}_n \log(5/\bar{X}_n)] \xleftarrow{\text{d}} \chi_1^2 \text{ under } H_0.$$

Simulation vs asymptotic approx

- ① Generate 10,000 samples, each sample has  $\sim 25$  iid  $\text{Poisson}(\lambda)$ .  
 $\text{Poisson}(5)$

- ② For each sample, we calculate the  $\bar{X}_n$  and  $-2 \log \lambda(\bar{X}_n)$ .



- ③ Plot the histogram using all the likelihood ratio values.

## Generalized LRT example

$$R = \{ \lambda(\bar{x}_n) \leq c \} = \{ -2\log \lambda(\bar{x}_n) \geq c' \}$$

$$p\text{-value} = P(R | H_0) = P(-2\log \lambda(\bar{x}_n) \geq c' | \lambda=5)$$

**Example 6.** Let  $X_1, \dots, X_{25}$  be i.i.d  $\text{Poisson}(\lambda)$ . Consider

$$H_0 : \lambda = 5 \leftrightarrow H_1 : \lambda \neq 5.$$

If we observe  $\bar{X}_{25} = 5.93$ , do we reject  $H_0$ ?

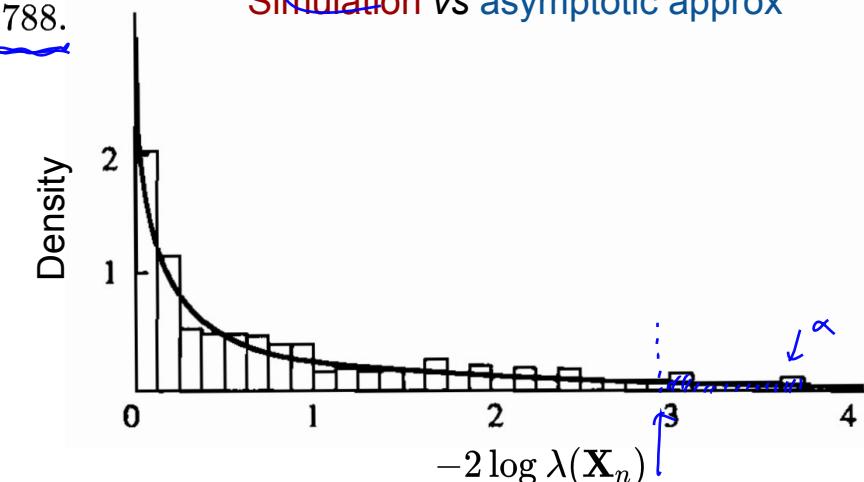
$$c' = qchisq(\alpha, df=1, \text{lower.tail}=\text{FALSE}) \Leftarrow \alpha = P(R | H_0) = pchisq(c', df=1, \text{lower.tail}=\text{FALSE})$$

$$-2\log \lambda(\bar{X}_n) = -50[(5.93 - 5) + 5.93 \cdot \log(5/5.93)] = 4.0788.$$

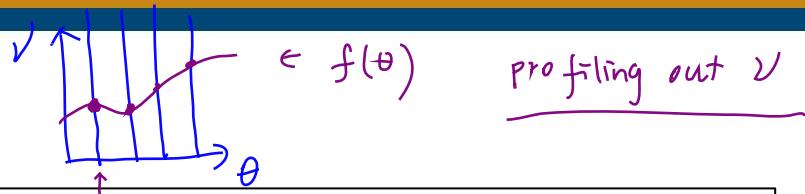
Simulation vs asymptotic approx

| Upper $\alpha$ quantile | 0.20  | 0.10  | 0.05  | 0.01  |
|-------------------------|-------|-------|-------|-------|
| Simulated               | 1.630 | 2.726 | 3.744 | 6.304 |
| $\chi^2$                | 1.642 | 2.706 | 3.841 | 6.635 |

$\uparrow$        $\uparrow$        $\uparrow$        $\uparrow$   
 $\{ -2\log \lambda(\bar{x}_n) \geq 1.642 \}$   
 Reject    Reject    Reject    Fail to Reject



## Generalized LRT example



**Example 7 (HW3 Prob 6).** Let  $X_1, \dots, X_n$  be i.i.d  $\text{Pareto}(\theta, \nu)$  with  $\nu$  being unknown. Consider

$$H_0 : \theta = 1 \leftrightarrow H_1 : \theta \neq 1.$$

Derive the expression for  $-2 \log \lambda(\mathbf{X}_n)$ .

$$f(x|\theta) = \begin{cases} \frac{\theta\nu^\theta}{x^{\theta+1}}, & \text{if } x \geq \nu, \\ 0, & \text{otherwise.} \end{cases}$$

Solution.  $L(\theta, \nu | \mathbf{x}_n) = \frac{\theta^n \nu^{n\theta}}{\left(\prod_{i=1}^n x_i\right)^{\theta+1}} \mathbb{1}\{X_{(1)} \geq \nu\}$

$$\Theta_0 = \left\{ \theta = 1, \nu > 0 \right\}, \quad \Theta = \left\{ \theta > 0, \nu > 0 \right\}$$

$$\sup_{\Theta_0} L(\theta, \nu | \mathbf{x}_n) = \sup_{\nu > 0} \frac{\nu^n}{\left(\prod_{i=1}^n x_i\right)^2} \mathbb{1}\{X_{(1)} \geq \nu\} = \frac{\nu^n}{\left(\prod_{i=1}^n x_i\right)^2} \quad \nu = X_{(1)}$$

$$= \frac{\sum \log(x_i/X_{(1)})}{n}$$

$$l(\theta, \nu | \mathbf{x}_n) = \begin{cases} n(\theta) + n\theta \log \nu - (\theta+1) \sum_{i=1}^n \log x_i & \text{if } \nu \leq X_{(1)}, \\ -\infty & \text{if } \nu > X_{(1)}. \end{cases} e^{\hat{\theta}_{MLE} \sum_{i=1}^n \log \frac{x_{(1)}}{x_i}}$$

$$\hat{\nu}_{MLE} = \frac{\sum \log x_i - \log X_{(1)}}{n}$$

$$\text{Fix } \theta, \max_{\nu \leq X_{(1)}} l(\theta, \nu | \mathbf{x}_n) = n \log \theta + n \theta \log X_{(1)} - (\theta+1) \sum_{i=1}^n \log x_i = f(\theta)$$

$$\max_{\theta > 0} f(\theta) = \sup_{\Theta} L(\theta, \nu | \mathbf{x}_n) = \left[ \frac{n}{\log \frac{\prod_{i=1}^n x_i}{X_{(1)}}} \right]^n \left[ \frac{X_{(1)}}{X_{(1)} - \sum_{i=1}^n \log x_i} \right] = \left[ \frac{\log \left( \frac{\prod_{i=1}^n x_i}{X_{(1)}} \right)}{n} \right]^n \in \frac{\partial f(\theta)}{\partial \theta} = \frac{n}{\theta} - \frac{n}{\sum_{i=1}^n \log x_i} \geq 0$$

# Generalized LRT example

$$\hat{\theta}_{MLE} = \frac{n}{\sum_{i=1}^n \log(x_i/x_{(1)})}$$

**Example 7 (HW3 Prob 6).** Let  $X_1, \dots, X_n$  be i.i.d  $Pareto(\theta, \nu)$  with  $\nu$  being unknown. Consider

$$H_0 : \theta = 1 \leftrightarrow H_1 : \theta \neq 1.$$

Derive the expression for  $-2 \log \lambda(\mathbf{X}_n)$ .  $\cancel{\downarrow X_1^2}$

Solution cont'd.

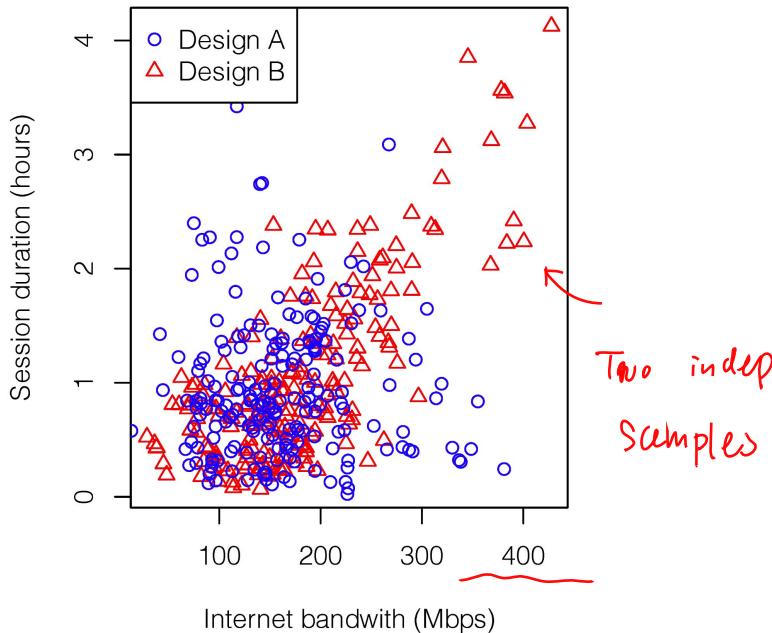
$$\begin{aligned} \sup_{\theta} L(\theta, \nu | \mathbf{x}_n) &= \left[ \frac{\log \prod_{i=1}^n \frac{x_i}{x_{(1)}}}{n} \right]^n e^{\hat{\theta}_{MLE} \sum_{i=1}^n \log(x_i/x_{(1)})} / \prod_{i=1}^n x_i \quad \text{Sampling distribution is difficult,} \\ \lambda(\mathbf{x}_n) &= \frac{x_{(1)}}{\left[ \frac{\log \prod_{i=1}^n \frac{x_i}{x_{(1)}}}{n} \right]^n e^n} / \left[ \frac{\log \prod_{i=1}^n \frac{x_i}{x_{(1)}}}{n} \right]^n e^n / \prod_{i=1}^n x_i \\ &= n^n e^{-n \sum_{i=1}^n \frac{x_{(1)}}{x_i}} \times \left[ \log \frac{\prod_{i=1}^n x_i}{\prod_{i=1}^n x_{(1)}} \right]^n \quad \text{Denote } T(\mathbf{x}_n) = \log \left( \frac{\prod_{i=1}^n x_i}{\prod_{i=1}^n x_{(1)}} \right) \\ &= n^n e^{-n} e^{-T(\mathbf{x}_n)} \times T(\mathbf{x}_n)^n \\ R = \{ \lambda(\mathbf{x}_n) \leq c \} &= \{ T(\mathbf{x}_n) \leq t_1 \text{ or } T(\mathbf{x}_n) \geq t_2 \} \quad \leftarrow f(t) = e^{-t} t^n \end{aligned}$$

# Two independent samples (*Post-midterm*)

*11.2 of Rice*

07/13/2021

# Youtube A/B testing

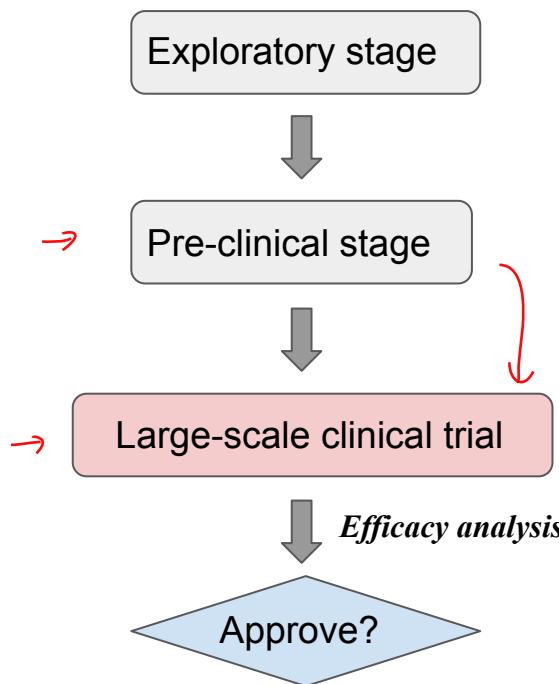


- Seems like Design B has more users with fast Internet.
- Was the A/B testing fully randomized?
- Denote mean bandwidth in Design A by  $\mu_A$ , and similarly denote  $\mu_B$ . Can we test

$$H_0 : \mu_A = \mu_B \text{ versus } H_1 : \mu_A < \mu_B ?$$



# Vaccine Testing and Approval



*Large, randomized, double-blind, placebo-controlled trial:*

|                     | Pfizer / BNT162b2 | Placebo |
|---------------------|-------------------|---------|
| SARS-CoV-2 infected | 9                 | 169     |
| No infection        | 21,711            | 21,559  |
| Total               | 21,720            | 21,728  |

\* Age  $\geq 16$ , infections observed with onset at least 7 days after the second dose;

\* Phase 3 data from [Pfizer.com](https://Pfizer.com)

$$\begin{aligned} H_0 : p_{\text{pfizer}} &= p_{\text{placebo}} & \text{versus} & \quad H_1 : p_{\text{pfizer}} \neq p_{\text{placebo}}. \\ H_0 : p_{\text{pfizer}} &\geq p_{\text{placebo}}/3 & \text{versus} & \quad H_1 : p_{\text{pfizer}} < p_{\text{placebo}}/3. \end{aligned}$$

$p_{\text{pfizer}} < p_{\text{placebo}}/3$

# Independent samples under two Normal populations

**Example 1.** Let  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} N(\mu_X, \sigma^2)$ ,  $Y_1, \dots, Y_m \stackrel{\text{iid}}{\sim} N(\mu_Y, \sigma^2)$  and the variance is unknown.  
 Consider the hypotheses

$$H_0 : \mu_X = \mu_Y \quad \text{versus} \quad H_1 : \mu_X \neq \mu_Y.$$

Derive the LRT statistic and its rejection region.

Solution. The joint likelihood is :  $L(\mu_X, \mu_Y, b^2 | \bar{X}_n, \bar{Y}_m) = \left( \frac{1}{\sqrt{2\pi b^2}} \right)^{m+n} e^{-\frac{\sum_{i=1}^n (x_i - \mu_X)^2}{2b^2} - \frac{\sum_{i=1}^m (y_i - \mu_Y)^2}{2b^2}}$

$$\textcircled{H}_0 = \{ \mu_X = \mu_Y, b^2 > 0 \}, \quad \textcircled{H} = \{ \mu_X \in \mathbb{R}, \mu_Y \in \mathbb{R}, b^2 > 0 \} \leftarrow$$

To calculate the unrestricted maximum likelihood  $\sup_{\textcircled{H}} L(\mu_X, \mu_Y, b^2 | \bar{X}_n, \bar{Y}_m)$ :

$$L(\mu_X, \mu_Y, b^2 | \bar{X}_n, \bar{Y}_m) = -\frac{m+n}{2} \log(2\pi b^2) - \frac{\sum_{i=1}^n (x_i - \mu_X)^2 + \sum_{i=1}^m (y_i - \mu_Y)^2}{2b^2}$$

$$\left\{ \begin{aligned} \frac{\partial L}{\partial \mu_X} &= \frac{\sum (x_i - \mu_X)}{b^2} = 0 \Rightarrow \hat{\mu}_X = \bar{X}_n \end{aligned} \right.$$

$$\left\{ \begin{aligned} \frac{\partial L}{\partial \mu_Y} &= \frac{\sum (y_i - \mu_Y)}{b^2} = 0 \Rightarrow \hat{\mu}_Y = \bar{Y}_m \\ \frac{\partial L}{\partial b^2} &= -\frac{m+n}{2} \cdot \frac{2}{b^2} + \frac{\sum (x_i - \mu_X)^2 + \sum (y_i - \mu_Y)^2}{b^3} \end{aligned} \right.$$

pooled variance estimate

$$\hat{b}^2_p = \frac{\sum (x_i - \bar{X}_n)^2 + \sum (y_i - \bar{Y}_m)^2}{m+n}$$

$$\sum_{i=1}^n (x_i - \bar{x})^2 + \sum_{i=1}^m (y_i - \bar{y})^2 = \sum_{i=1}^n (x_i - \bar{x}_n)^2 + n(\bar{x}_n - \bar{x})^2 + \sum_{i=1}^m (y_i - \bar{y}_m)^2 + m(\bar{y}_m - \bar{y})^2$$

$\downarrow x_i \rightarrow \bar{x}_n + \bar{x}_n - \bar{x}$

$$= \left( \frac{n+m}{m+n} \right) \hat{\sigma}_p^2 + \left( \frac{m}{m+n} \right) (\bar{x}_n - \bar{x})^2 + \left( \frac{n}{m+n} \right) (\bar{y}_m - \bar{y})^2$$

## Independent samples under two Normal populations

Solution cont'd.

$$\begin{aligned} \sup_{\mathbb{H}} L(\mu_x, \mu_y, \sigma^2 | \bar{x}_n, \bar{y}_m) &= L(\hat{\mu}_x, \hat{\mu}_y, \hat{\sigma}_p^2 | \bar{x}_n, \bar{y}_m) \\ &= \left( \frac{1}{\sqrt{2\pi}\hat{\sigma}_p} \right)^{m+n} e^{-\frac{\sum(x_i - \bar{x}_n)^2 + \sum(y_i - \bar{y}_m)^2}{2\hat{\sigma}_p^2}} = \left( \frac{1}{\sqrt{2\pi}\hat{\sigma}_p} \right)^{m+n} e^{-\frac{m+n}{2}} \end{aligned}$$

To calculate the restricted maximum likelihood  $\sup_{\mathbb{H}_0} L(\mu_x, \mu_y, \sigma^2 | \bar{x}_n, \bar{y}_m)$ :

$$H_0: \mu_x = \mu_y, \quad x_1, \dots, x_n \sim N(\mu_x, \sigma^2), \quad y_1, \dots, y_m \sim N(\mu_x, \sigma^2)$$

$$\begin{aligned} \sup_{\mathbb{H}_0} L(\mu_x, \mu_y, \sigma^2 | \bar{x}_n, \bar{y}_m) &= L(\hat{\mu}, \hat{\mu}, \hat{\sigma}^2) \\ &= \left( \frac{1}{\sqrt{2\pi}\hat{\sigma}} \right)^{m+n} e^{-\frac{\sum(x_i - \hat{\mu})^2 + \sum(y_i - \hat{\mu})^2}{2\hat{\sigma}^2}} = \left( \frac{1}{\sqrt{2\pi}\hat{\sigma}} \right)^{m+n} e^{-\frac{m+n}{2}} \end{aligned}$$

$$\lambda(\bar{x}_n, \bar{y}_m) = \left( -\frac{\hat{\sigma}_p^2}{\hat{\sigma}^2} \right)^{\frac{m+n}{2}} = \left( \frac{\frac{\sum(x_i - \bar{x}_n)^2 + \sum(y_i - \bar{y}_m)^2}{m+n}}{\frac{\sum(x_i - \hat{\mu})^2 + \sum(y_i - \hat{\mu})^2}{m+n}} \right)^{\frac{m+n}{2}}$$

From the previous page,

$$\begin{aligned}\sum (x_i - \bar{x})^2 + \sum (f_i - \bar{f}_n)^2 &= (m+n) \hat{b}_p^2 + n(\bar{x}_n - \bar{\mu})^2 + m(\bar{f}_m - \bar{\mu})^2 \\&\quad n \left( \bar{x}_n - \frac{n\bar{x}_n + m\bar{f}_m}{n+m} \right)^2 \\&\quad m \left[ \frac{n(\bar{x}_n - \bar{f}_m)}{n+m} \right]^2 \\&= (m+n) \hat{b}_p^2 + \frac{mn}{(n+m)^2} (\bar{x}_n - \bar{f}_m)^2\end{aligned}$$

$$\lambda(\bar{x}_n, \bar{f}_m) = \left[ \frac{(m+n) \hat{b}_p^2}{(m+n) \hat{b}_p^2 + \frac{mn}{m+n} (\bar{x}_n - \bar{f}_m)^2} \right]^{\frac{m+n}{2}} = \left[ \frac{1}{1 + \frac{mn}{(m+n)^2} \frac{(\bar{x}_n - \bar{f}_m)^2}{\hat{b}_p^2}} \right]^{\frac{m+n}{2}}$$

$$R = \{ \lambda(\bar{x}_n, \bar{f}_m) \leq c \} = \left\{ \frac{mn}{(m+n)^2} \frac{(\bar{x}_n - \bar{f}_m)^2}{\hat{b}_p^2} \geq c \right\}$$

$$\left\{ \frac{|\bar{x}_n - \bar{f}_m|}{S_p \sqrt{\frac{1}{n} + \frac{1}{m}}} \geq \sqrt{c} \right\} = \left\{ \sqrt{\frac{mn}{m+n}} \cdot \frac{|\bar{x}_n - \bar{f}_m|}{\sqrt{\sum (x_i - \bar{x}_n)^2 + \sum (f_i - \bar{f}_m)^2}} \geq \sqrt{c} \right\} \xrightarrow{(n+m \rightarrow)} S_p \xrightarrow{(n+m \rightarrow)} \sqrt{c}$$

consistent with our intuition.

# Independent samples under two Normal populations

**Theorem A.** Let  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} N(\mu_X, \sigma^2)$ ,  $Y_1, \dots, Y_m \stackrel{\text{iid}}{\sim} N(\mu_Y, \sigma^2)$  and the variance is unknown.

Then

$$\sum_{i=1}^n (X_i - \bar{X}_n)^2 + \sum_{i=1}^m (Y_i - \bar{Y}_m)^2 \sim \sigma^2 \chi_{m+n-2}^2,$$

and it is mutually independent from  $\bar{X}_n$  and  $\bar{Y}_m$ .

Proof.  $\{X_1, \dots, X_n\} \stackrel{\text{iid}}{\sim} N(\mu_X, \sigma^2)$

$$\frac{1}{b^2} \sum_{i=1}^n (X_i - \bar{X}_n)^2 \stackrel{\text{indep from } \bar{X}_n}{\sim} \chi_{n-1}^2$$

$\{Y_1, \dots, Y_m\} \stackrel{\text{iid}}{\sim} N(\mu_Y, \sigma^2)$

$$\frac{1}{b^2} \sum_{i=1}^m (Y_i - \bar{Y}_m)^2 \stackrel{\text{indep from } \bar{Y}_m}{\sim} \chi_{m-1}^2$$

$$\frac{\sum (X_i - \bar{X}_n)^2 + \sum (Y_i - \bar{Y}_m)^2}{b^2} \sim \chi_{n-1}^2 + \chi_{m-1}^2 = \chi_{n+m-2}^2$$

$$\perp \bar{X}_n \perp \bar{Y}_m.$$

$$\hat{\sigma}_n^2, \hat{s}_n^2$$

$$Z \sim N(0,1)$$

$$\frac{V}{2} \sim \chi_{k^2}^2$$

Independent samples under two Normal populations

$$\frac{\sqrt{V}}{\sqrt{k}} \sim t_k,$$

**Theorem B.** Let  $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu_X, \sigma^2)$ ,  $Y_1, \dots, Y_m \stackrel{iid}{\sim} N(\mu_Y, \sigma^2)$  and the variance is unknown.

Then

$$\hat{\sigma}_p^2 = \frac{\sum_{i=1}^{n+m} (X_i - \bar{X}_n)^2 + (Y_i - \bar{Y}_m)^2}{n+m}$$

$$\frac{(\bar{X}_n - \bar{Y}_m) - (\mu_X - \mu_Y)}{S_p \sqrt{\frac{1}{n} + \frac{1}{m}}} \sim t_{n+m-2}.$$

$$\rightarrow S_p^2 = \frac{\sum_{i=1}^n (X_i - \bar{X}_n)^2 + \sum_{i=1}^m (Y_i - \bar{Y}_m)^2}{n+m-2}$$

$$\bar{X}_n \sim N(\mu_X, \frac{\sigma^2}{n}), \quad \bar{Y}_m \sim N(\mu_Y, \frac{\sigma^2}{m})$$

$$\Rightarrow \bar{X}_n - \bar{Y}_m \sim N(\mu_X - \mu_Y, \sigma^2(\frac{1}{n} + \frac{1}{m}))$$

$$\frac{1}{\sigma^2} \left\{ \sum_{i=1}^n (X_i - \bar{X}_n)^2 + \sum_{i=1}^m (Y_i - \bar{Y}_m)^2 \right\} \sim \chi^2_{n+m-2}$$

$$\frac{(\bar{X}_n - \bar{Y}_m) - (\mu_X - \mu_Y)}{\sigma \sqrt{\frac{1}{n} + \frac{1}{m}}} \sim N(0,1)$$

$$= \frac{(n+m-2) S_p^2}{\sigma^2} \sim \chi^2_{n+m-2}$$

$$\frac{(\bar{X}_n - \bar{Y}_m) - (\mu_X - \mu_Y)}{\sigma \sqrt{\frac{1}{n} + \frac{1}{m}}} \sim N(0,1)$$

$$= \frac{\sqrt{\frac{(n+m-2) S_p^2}{\sigma^2}} / \sqrt{n+m-2}}{\sigma \sqrt{\frac{1}{n} + \frac{1}{m}}}$$

$$\sim t_{n+m-2}$$

$$= \frac{(\bar{X}_n - \bar{Y}_m) - (\mu_X - \mu_Y)}{S_p \sqrt{\frac{1}{n} + \frac{1}{m}}}$$

# Independent samples under two Normal populations

**Example 1 cont'd.** Let  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} N(\mu_X, \sigma^2)$ ,  $Y_1, \dots, Y_m \stackrel{\text{iid}}{\sim} N(\mu_Y, \sigma^2)$  and the variance is unknown. Consider the hypotheses

$$H_0 : \mu_X = \mu_Y \quad \text{versus} \quad H_1 : \mu_X \neq \mu_Y. \quad \frac{|\bar{X}_n - \bar{Y}_m|}{S_p \sqrt{\frac{1}{n} + \frac{1}{m}}}$$

Derive the test with significance level  $\alpha$ .

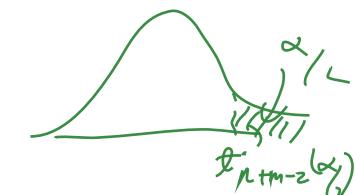
Solution. We derived the rejection region:  $R = \left\{ \sqrt{\frac{mn}{m+n}} \frac{|\bar{X}_n - \bar{Y}_m|}{\sqrt{\sum_{i=1}^n (X_i - \bar{X}_n)^2 + \sum_{i=1}^m (Y_i - \bar{Y}_m)^2}} \geq c'' \right\}$ .

Under  $H_0$ ,  $\mu_X - \mu_Y = 0$  and  $\frac{\bar{X}_n - \bar{Y}_m}{S_p \sqrt{\frac{1}{n} + \frac{1}{m}}} \sim t_{n+m-2}$

To get the significance level  $\alpha$ :

$$\alpha = P(R | H_0) = P\left(\frac{|\bar{X}_n - \bar{Y}_m|}{S_p \sqrt{\frac{1}{n} + \frac{1}{m}}} \geq c'' \mid \mu_X = \mu_Y\right)$$

$\Rightarrow c'' = t_{n+m-2}(\alpha)$ ,



# Independent samples under two Normal populations

**Corollary B.** Let  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} N(\mu_X, \sigma^2)$ ,  $Y_1, \dots, Y_m \stackrel{\text{iid}}{\sim} N(\mu_Y, \sigma^2)$  and the variance is unknown.

Then  $(1 - \alpha) \times 100\%$  exact CI for  $\mu_X - \mu_Y$  is

$$(\bar{X}_n - \bar{Y}_m) \pm t_{n+m-2}(\alpha/2) \cdot S_p \sqrt{\frac{1}{n} + \frac{1}{m}}.$$

$$\begin{aligned} & \frac{(\bar{X}_n - \bar{Y}_m) - (\mu_X - \mu_Y)}{S_p \sqrt{\frac{1}{n} + \frac{1}{m}}} \sim t_{n+m-2} \\ & 1 - \alpha = P \left( \left| \frac{(\bar{X}_n - \bar{Y}_m) - (\mu_X - \mu_Y)}{S_p \sqrt{\frac{1}{n} + \frac{1}{m}}} \right| \leq t_{n+m-2}(\alpha/2) \right) \\ & = P \left[ (\bar{X}_n - \bar{Y}_m) - t_{n+m-2}(\alpha/2) \cdot S_p \sqrt{\frac{1}{n} + \frac{1}{m}} \leq \mu_X - \mu_Y \leq (\bar{X}_n - \bar{Y}_m) + t_{n+m-2}(\alpha/2) \cdot S_p \sqrt{\frac{1}{n} + \frac{1}{m}} \right] \end{aligned}$$

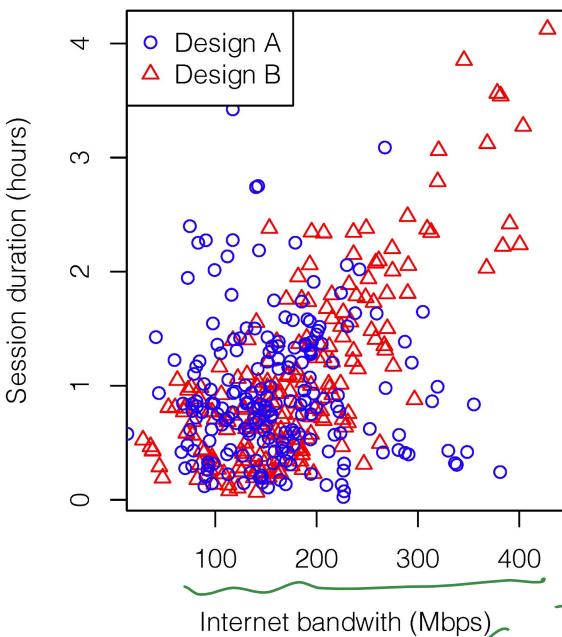
# Independent samples under two Normal populations

**Example 1 cont'd.** Let  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} N(\mu_X, \sigma^2)$ ,  $Y_1, \dots, Y_m \stackrel{\text{iid}}{\sim} N(\mu_Y, \sigma^2)$  and the variance is unknown. The LRTs with significance level  $\alpha$  for the following hypotheses can be derived.

$$\begin{aligned}
 H_0: \mu_X = \mu_Y &\Leftrightarrow H_1: \mu_X \neq \mu_Y. \Leftrightarrow R = \left\{ \frac{|\bar{X}_n - \bar{Y}_m|}{S_p \sqrt{\frac{1}{n} + \frac{1}{m}}} \geq t_{n+m-2}(\alpha/2) \right\} \\
 A = \left\{ (\bar{X}_n - \bar{Y}_m) - t_{n+m-2}(\alpha/2) \cdot S_p \sqrt{\frac{1}{n} + \frac{1}{m}} \leq \mu_X - \mu_Y \leq (\bar{X}_n - \bar{Y}_m) + t_{n+m-2}(\alpha/2) \cdot S_p \sqrt{\frac{1}{n} + \frac{1}{m}} \right\} \\
 &= \left\{ \frac{|\bar{X}_n - \bar{Y}_m|}{S_p \sqrt{\frac{1}{n} + \frac{1}{m}}} \leq t_{n+m-2}(\alpha/2) \right\} \\
 H_0: \mu_X = \mu_Y &\Leftrightarrow H_1: \mu_X > \mu_Y. \Leftrightarrow R = \left\{ \frac{\bar{X}_n - \bar{Y}_m}{S_p \sqrt{\frac{1}{n} + \frac{1}{m}}} \geq t_{n+m-2}(\alpha) \right\}
 \end{aligned}$$

$$H_0: \mu_X = \mu_Y \Leftrightarrow H_1: \mu_X < \mu_Y. \Leftrightarrow R = \left\{ \frac{\bar{X}_n - \bar{Y}_m}{S_p \sqrt{\frac{1}{n} + \frac{1}{m}}} \leq -t_{n+m-2}(\alpha) \right\}$$

# Youtube A/B testing



$$n = 200, m = 200, n+m-2 = 398.$$

$$t_{398}$$

$$p\text{-value} = P\left(\frac{\bar{X}_{200} - \bar{Y}_{200}}{S_p \sqrt{\frac{1}{200} + \frac{1}{200}}} \leq -2.160845 | H_0\right) = pt(-2.160845, df=398)$$

$$H_0: \mu_A = \mu_B \text{ versus } H_1: \mu_A < \mu_B ?$$

$$\alpha = 0.05$$

$$= 0.0157 < \alpha.$$

Assume  $X_1, \dots, X_{200} \sim N(\mu_A, \sigma^2), Y_1, \dots, Y_{200} \sim N(\mu_B, \sigma^2)$ .

The LRT rejection region is  $R = \left\{ \frac{\bar{X}_{200} - \bar{Y}_{200}}{S_p \sqrt{\frac{1}{200} + \frac{1}{200}}} \leq -t_{398}(\alpha) \right\}$ .

```
> mean(A_bandwidth)
[1] 3.108194
> mean(B_bandwidth)
[1] 3.384364
> {sum((A_bandwidth-mean(A_bandwidth))^2) + sum((B_bandwidth-mean(B_bandwidth))^2)}/398
[1] 1.633449  $t_{398}(\alpha)$ 
> qt(0.05, df=398, lower.tail=FALSE)
[1] 1.648691
```

$$\frac{\bar{X}_{200} - \bar{Y}_{200}}{S_p \sqrt{\frac{1}{200} + \frac{1}{200}}} = \frac{3.10894 - 3.384364}{\sqrt{1.633449} \sqrt{\frac{1}{100}}} = -2.160845$$

Conclusion: We reject  $H_0$ , and conclude that there is enough evidence to suggest that  $\mu_A < \mu_B$ .



## Byzantine church wood

$$P\text{-value} = P\left(\frac{|\bar{X}_n - \bar{Y}_m|}{S_p \sqrt{\frac{1}{n} + \frac{1}{m}}} \geq 1.2899 \mid H_0\right) = 2 * pt(-1.2899, df=2) = 0.21109$$

$\alpha = 0.05$

**Example 2.** Samples of wood were obtained from the core and periphery of a Byzantine church. The date of the wood were determined.

$$\bar{X}_n = 1249.86 \quad \bar{Y}_m = 1261.33$$

$$S_p^2 = 433.13$$

| Core | Periphery |
|------|-----------|
| 1294 | 1251      |
| 1279 | 1248      |
| 1274 | 1240      |
| 1264 | 1232      |
| 1263 | 1220      |
| 1254 | 1218      |
| 1251 | 1210      |

$$\begin{matrix} \uparrow \\ n=14 \end{matrix}$$

$$H_0 : \mu_C = \mu_P \text{ versus } H_1 : \mu_C \neq \mu_P ?$$

The rejection region is  $R = \left\{ \frac{|\bar{X}_n - \bar{Y}_m|}{S_p \sqrt{\frac{1}{n} + \frac{1}{m}}} \geq t_{n+m-2}(\alpha/2) \right\}$ .

$$t_{n+m-2(\alpha/2)} = t_{21}(0.05/2) = 2.0796.$$

$$\frac{|\bar{X}_n - \bar{Y}_m|}{S_p \sqrt{\frac{1}{n} + \frac{1}{m}}} = \frac{|1249.86 - 1261.33|}{\sqrt{433.13} \sqrt{\frac{1}{14} + \frac{1}{9}}} = 1.2899.$$

$n=9$  we fail to reject  $H_0$ , and conclude that there is not enough evidence to support the alternative hypothesis.

Tomorrow ...

- Non-parametric approach - Mann-Whitney test
- Comparing **paired samples**

