

## Recitation Two

### PART ONE: Testing

#### Basics of z-tests

We have a  $z$ -stat that we believe is approximately normal under the null. For now, I will assume if  $\tau = \frac{\hat{\theta} - \theta_0}{SE(\hat{\theta})} \xrightarrow{d} N(0, 1)$  if  $\theta_0 = E[\hat{\theta}]$  truly.

More specifically, say  $\hat{\theta} = \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$

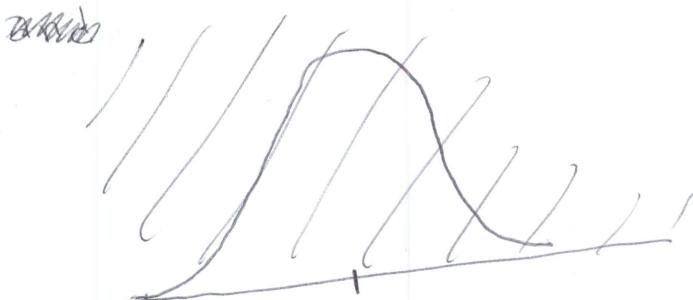
$$\tau = \frac{\bar{x} - \mu_0}{SE(\bar{x})} \xrightarrow{d} N(0, 1)$$

OR

$$\frac{\bar{x} - \mu_0}{SE(\bar{x})/\sqrt{n}} \xrightarrow{d} N(0, 1) \text{ if } n \text{ large.}$$

Then under the assumption that

$$E[x_i] = \mu_0, \quad \frac{\bar{x} - \mu_0}{SE(\bar{x})/\sqrt{n}} \sim N(0, 1)$$



$$\text{or } \bar{x} \sim N(\mu_0, \frac{Var(x_i)}{n})$$

Distribution of  $\bar{x}$  assuming null is true.

Variance of the distribution depends on  $Var(x_i)$  and  $n$ .



Our question is:

2

Is the null we assumed reasonable?

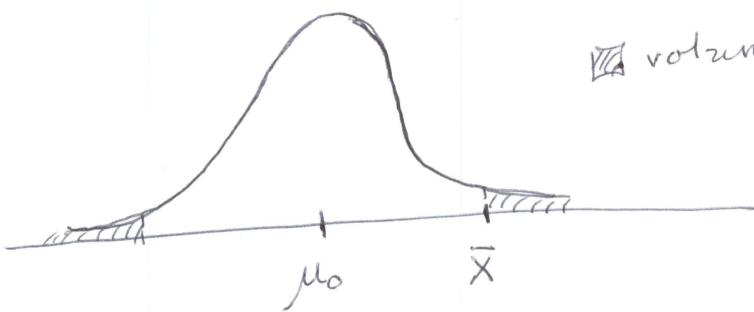
Another way of asking this is:

Assuming the null is true, how likely would our data we observe truly be?

How do we answer this? p-value!

p-value is the probability of seeing something more extreme than the same with the data assuming the null is correct.

$$\text{volume} = \text{p-value}.$$



So, we think: if the prob. of seeing a more extreme value is low, then the null is ~~not~~ very reasonable.

The critical value is simply a value that gives us a specified p-value. This specified p-value is called a significance level)

Reject null if p-value < sig. level  
OR  $|x| < \text{crit. value}$ .

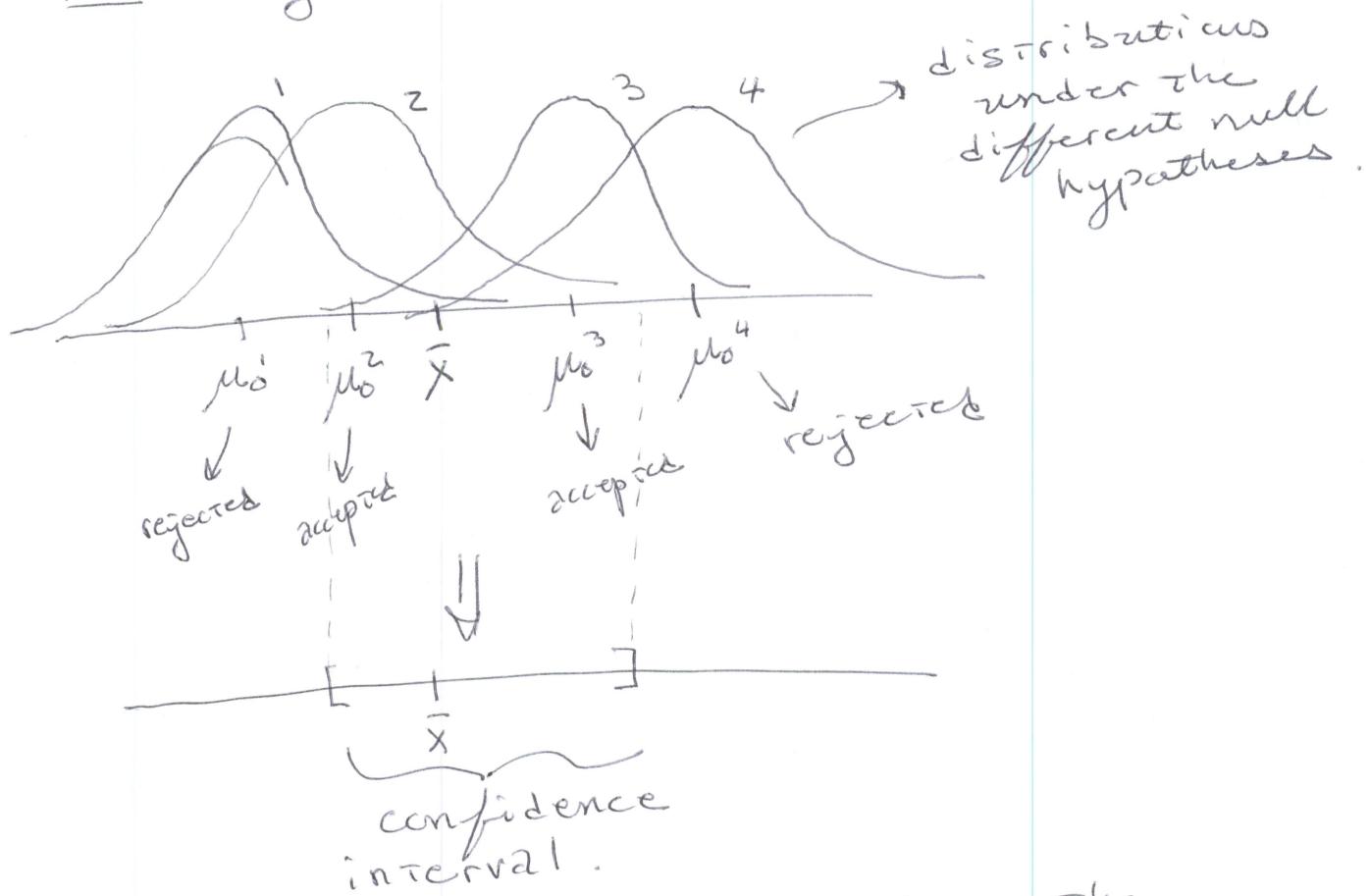
These are the same thing.

## Finally, CONFIDENCE INTERVALS

3

A CONFIDENCE INTERVAL is an interval that we are  $(1-\alpha)\%$  confident that the true parameter value is in that interval. ( $\alpha$  is the confidence level - you choose this - usually  $\alpha = .05$ )

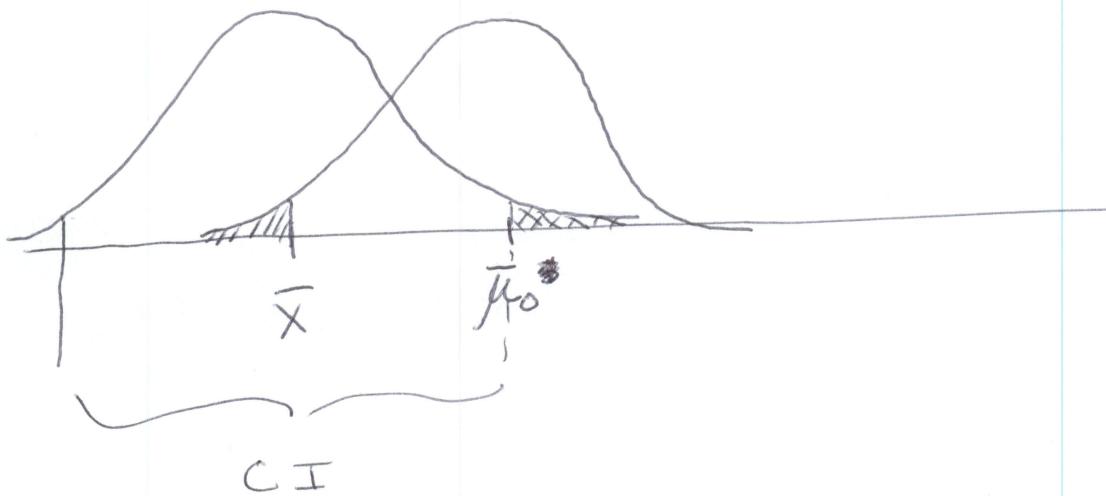
For a ~~test~~ sample mean test, one such confidence interval is the set of null hypothesis values  $\mu_0$  such that  $H_0: \mu = \mu_0$  is not rejected.



It just so happens that, since the normal distribution is symmetric, this set can be calculated as

$$\bar{X} \pm SE(\bar{X}) Z_{1-\alpha/2}$$

Why?



All  $\mu_0 > \bar{\mu}_0$  are rejected

if and only if

All  $\mu_0 > \bar{\mu}_0$  are greater than

$$\bar{x} + SE(\bar{x}) Z_{1-\alpha/2}$$

If  $\bar{\mu}_0$  is the upper bound on the confidence interval, both the area

of  $\blacksquare$  and  $\blacksquare$  are equal to  $\frac{\alpha}{2}$ .

If they got further apart, both would be less than  $\frac{\alpha}{2}$ .

Now, going back to the beginning, this can all be re-analyzed in the exact same fashion, but instead of setting  $\hat{\theta} = \bar{x}$ , we set  $\hat{\theta} = \hat{\beta}_1$  (from OLS) and instead of  $\hat{\theta} = \bar{x}$ . And we set  $SE(\hat{\theta}) = SE(\hat{\beta}_1)$

instead of  $SE(\hat{\beta}) = SE(\bar{x})$

One of the major points of the past few weeks of lectures is that:

$$1) \frac{\hat{\beta}_1 - \beta_1^*}{SE(\hat{\beta}_1)} \xrightarrow{d} N(0, 1)$$

that is, we can use the analysis talked about before

$$2) SE(\hat{\beta}_1) = \sqrt{\frac{1}{n} \frac{\text{var}[(x_i - \mu_x) u_i]}{(\sigma_x^2)^2}}$$

where  $\mu_x$  is the TRUE mean of  $x_i$   
and  $\sigma_x^2$  is the TRUE var of  $x_i$

(when I say of  $x_i$ , I mean  
the distribution of  $x_i$ )

## TOPIC TWO : REGRESSION

6

OSS Assumptions: Without these, we cannot ensure  $\hat{\beta}_i$  is unbiased, consistent, asymptotically normally distributed, or anything else we would like. In short, without these assumptions, it is very hard to interpret  $\hat{\beta}_i$ .

$$1) \mathbb{E}[u_i | X_i = x] = 0$$

Means  $u_i$  and  $X_i$  do not correlate.  
 $X_i$  cannot "explain"  $u_i$  at all.

$$2) (X_i, Y_i) \text{ are iid.}$$

That is person  $i$  is chosen randomly from the population. e.g.  $i$  and another draw are not siblings because we have family data.

$$3) \text{Large outliers are rare:}$$

$$\mathbb{E}[X_i^4] < \infty, \mathbb{E}[Y_i^4] < \infty$$

$$R^2: \hat{Var}(Y_i) = \hat{Var}(\hat{\beta}_0 + \hat{\beta}_1 X_i + \hat{u}_i)$$

$$\hat{u}_i = Y_i - \hat{\beta}_0 - \hat{\beta}_1 X_i \quad \hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$$

$$= \hat{Var}(\hat{Y}_i) + \hat{Var}(\hat{u}_i) + 2 \operatorname{cov}(\hat{Y}_i, \hat{u}_i)$$

$$= \frac{1}{n} \sum_{i=1}^n \hat{u}_i^2$$

$$TSS = ESS + RSS$$

so,  $R^2 = \frac{ESS}{TSS}$ , the proportion of variance in  $Y_i$  "explained" by  $X_i$

Final Note:  $SER = \frac{RSS}{n-k-1} = \text{estimate of var}(u_i)$  ( $n=1$  for us, so far)