

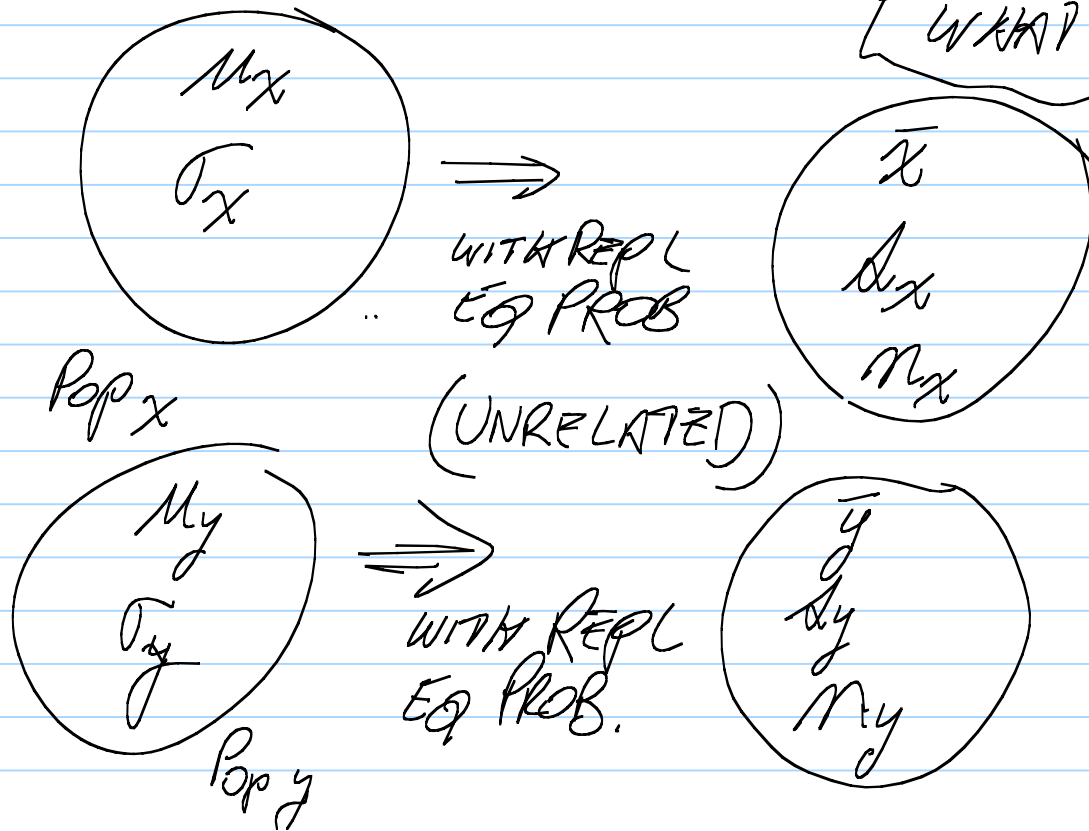
STAT 200 10-21-09

Note Title

10/21/2009

TODAY CI FOR DIFF μ_x, μ_y

IDEA



WHAT IS $\mu_x - \mu_y$

PETER DONNELLY
(TED SERIES)
GENOME
LETTERS A-ZTC

COIN TOSSES
HTH
HTT

QUESTION: WHICH TAKES LONGER ON AVG.

(a) WAIT FOR HTH

(b) WAIT FOR HTT

(c) SAME.

LET $X =$ # TOSSES TO GET
HTH

$Y =$ # TOSSES TO GET
HTT.

SAMPLES (STUDENT GENERATED)

x_1, x_2, \dots, x_{n_x} form \bar{x}, σ_x, n_x

y_1, y_2, \dots, y_{n_y} " \bar{y}, σ_y, n_y

95% CI FOR $\mu_x - \mu_y$:

$$\bar{x} - \bar{y} \pm 1.96 \sqrt{\sigma_x^2/n_x + \sigma_y^2/n_y}$$

{ CAUTION: SEPARATE CI
FOR X, Y WON'T
DO

EASIER TO GET
HTT THAN HTH
BECAUSE BOTH REQUIRE
LEADING OFF WITH
HT, BUT IF FAIL TO
GET HTT (GETTING HTH) YOU ARE ALREADY AHEAD

CLAIM: (a)
HT(T) ~~H~~ FOR H NEXT
WAITING FOR HTH?
HT(H)
WAITING FOR HTT?

(a) AVG WAIT TO HTA
IS LARGER ($M_x > M_y$)
(b) AVG WAIT TO HTT
IS LARGER ($M_x < M_y$)
(c) SAME AVG WAIT

CAN BE PROVEN THAT ON AVG TAKES TEN
TOSSES TO GET HTH BUT ONLY 8 TO GET HTT;

HHTTTHTTTTTHHTTTT(HTT)T
HTTTHHHHTHTH

MY SCORE $X = 32$

2) HTA 15 5 8 15 11 7 8 11 $\{x_i\}$ $\bar{x} = 10$
 $s_x = 3.66$
 $n_x = 8$

3) HTT 3 4 5 14 9 10 8 16 3 $\bar{y} = 8$
 $s_y = 4.74$
 $n_y = 9$

95% CI FOR μ_x $10 \pm 1.96 \frac{3.66}{\sqrt{8}}$

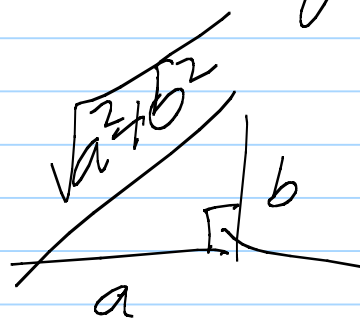
95% CI FOR μ_y $8 \pm 1.96 \frac{4.74}{\sqrt{9}}$

95% CI FOR $\mu_x - \mu_y = (10 - 8) \pm 1.96 \sqrt{\frac{3.66^2}{8} + \frac{4.74^2}{9}}$

$\bar{x} - \bar{y}$

$\sqrt{\frac{s_x^2}{n_x} + \frac{s_y^2}{n_y}}$

PYTHAGORAS' THEOREM



STATISTICAL INDEPENDENCE

CAVEATS:

REQUIRE n_x, n_y "LARGE ENOUGH"
[LATER WE DISCUSS "BOOTSTRAP" RESAMPLING YOUR DATA
TO CHECK UP ON STABILITY OF CI]

RETURN TO OUR CI FOR $\mu_x - \mu_y$:

$$(10-8) \pm 1.96 \sqrt{\frac{3.65^2}{8} + \frac{4.74^2}{9}}$$

{ -2, 6 } WAS (BEFORE YOU SEE IT)

$$P(\mu_x - \mu_y \text{ "IN" (COVERED BY) } (\bar{x} - \bar{y}) \pm 1.96 \sqrt{\frac{\sigma_x^2}{n_x} + \frac{\sigma_y^2}{n_y}}) \approx 0.95$$

TO GET MORE PRECISION (NARROWER CI)

REQUIRES n_x, n_y LARGER.

TO MAKE $\sqrt{\frac{\sigma_x^2}{n_x} + \frac{\sigma_y^2}{n_y}}$ SMALL

SHOULD ALLOCATE MORE DATA TO

SAMPLING x IF σ_x^2 IS LARGEST (SEE BELOW)

WHY THE $\frac{\sigma_x^2}{n_x} + \frac{\sigma_y^2}{n_y}$?

TO PREPARE YOU FOR THIS I'LL SPEAK TO
AN ARITHMETICAL FACT THAT UNDERLIES IT

LOOK AT TWO LISTS OF SCORES x, y .

x	y	$x - y$
7	4 6 9	$7-4=-3$ $7-6=-1$ $7-9=-2$
5		$5-4=1$ $5-6=-1$ $5-9=-4$
3		$3-4=-1$ $3-6=-3$ $3-9=-6$

$$d_{(x-y)} = \{ -3, -5, -8, 1, -1, -4, -1, -3, -6 \}$$

$$\text{CLAIM: } d_{(x-y)} = \sqrt{d_x^2 + d_y^2}$$

ROOTS OF THIS ARE IN FACT
 THAT FOR x "INDEPENDENT OF y " (ALL POSS COMBOS
 AS ABOVE)

WE HAVE $\overline{xy} = \bar{x}\bar{y}$

PROOF:

\equiv

$$\overline{xy} = \frac{\sum_i \sum_j x_i y_j}{n^2}$$

OCCUR IN ALL POSS
 COMBOS

$$= \frac{\sum_i x_i \sum_j y_j}{n \cdot n} = \bar{y} \frac{\sum x_i}{n} = \bar{y} \bar{x}$$

FURTHERMORE, LET'S JUST SUPPOSE $\bar{x} = \bar{y} = 0$

THEN $(x-y)^2 = x^2 - 2xy + y^2$

$$= \bar{x}^2 - 2\bar{x}\bar{y} \oplus \bar{y}^2 = \bar{x}^2 \oplus \bar{y}^2$$

$$\Rightarrow \overbrace{(x-y) - (\bar{x}-\bar{y})}^0}{}^2 = \overbrace{(x-\bar{x}) \oplus (y-\bar{y})}^2$$

(n-DIVISOR) $\text{Var}(x-y) = \text{Var}x \oplus \text{Var}y$

CASE OF
 x INDEP
 OF y (MEANS ALL COMBOS)

$$\Rightarrow \sigma_{(x-y)}^2 = \sigma_x^2 \oplus \sigma_y^2$$

SAMPLE ALLOCATION.

95% CI FOR $\mu_x - \mu_y$ IS $(\bar{x} - \bar{y}) \pm 1.96 \sqrt{\frac{\sigma_x^2}{n_x} + \frac{\sigma_y^2}{n_y}}$

[IF TOTAL SAMPLE SIZE $n = n_x + n_y$ IS FIXED

HOW SHOULD WE ALLOCATE TO MAKE CI NARROWEST?]

SEEMS WE SHOULD MAKE n_x LARGER IF $\sigma_x^2 \gg \sigma_y^2$.

FOR EXAMPLE. THE ACTUAL PRINCIPLE OF ALLOCATION

IS A BIT SUBTLE. LET $f = \frac{n_x}{n_x + n_y}$ BE THE

FRACTION OF SAMPLES DEVOTED TO X DATA.

HOW TO CHOOSE f TO MINIMIZE $\frac{\sigma_x^2}{n_x} \oplus \frac{\sigma_y^2}{n_y}$
THAT IS, CHOOSE f TO MINIMIZE

$$\frac{\sigma_x^2}{f n} \oplus \frac{\sigma_y^2}{(1-f)n} \quad \text{WHERE } n = n_x + n_y \text{ IS TOTAL SAMPLE SIZE.}$$

THE SOLUTION (REQUIRES CALCULUS) IS

$$f_x = \frac{\sigma_x}{\sigma_x + \sigma_y} = \text{IDEAL FRACTION OF } n \text{ DEVOTED TO SAMPLING } x \text{ DATA.}$$

OF COURSE WE DON'T KNOW σ_x, σ_y UNTIL AFTER