# Notes 10.2a:  Comparing Two Means

1. **If we want to compare the mean of some quantitative variable for the individuals in Population 1 and Population 2?**

   - The best approach is to take separate random samples from each population and to compare the sample means.

   - Our parameters of interest are the population means $\mu_1$ and $\mu_2$.

2. **Suppose we want to compare the average effectiveness of two treatments in a completely randomized experiment.**

   - The parameters $\mu 1$ and $\mu 2$ are the true mean responses for Treatment 1 and Treatment 2.  We use the mean response in the two groups to make the comparison.

   Here's a table that summarizes these two situations:

   | Population or treatment | Parameter | Statistic | Sample size |
   |---|---|---|---|
   | 1 | $\mu_1$ | $\bar{x}_1$ | $n_1$ |
   | 2 | $\mu_2$ | $\bar{x}_2$ | $n_2$ |

---

**The Sampling Distribution of the Difference Between Sample Means**

Choose
- an SRS of size $n_1$ from Population 1 with mean $\mu_1$ and std. dev. $\sigma_1$
- an SRS of size $n_2$ from Population 2 with mean $\mu_2$ and std. dev. $\sigma_2$.
- **The samples MUST be independent**

*ADDITIONAL CONDITION FOR 2 SAMPLE CI's & TESTS.*

**Center** The mean of the sampling distribution is an **unbiased estimator** of the difference in population means    $\mu_{\bar{x}_1-\bar{x}_2} = \mu_{\bar{x}_1} - \mu_{\bar{x}_2} = \mu_1 - \mu_2$
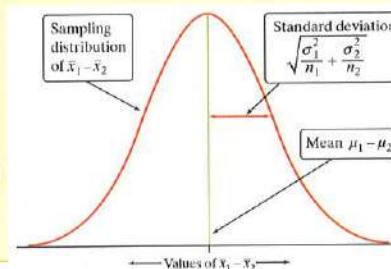
**Spread**    $\sigma_{\bar{x}_1-\bar{x}_2} = \sqrt{\dfrac{\sigma_1^2}{n_1} + \dfrac{\sigma_1^2}{n_2}}$

*This formula on Green sheet!*

As long as *each sample* is no more than 10% of its population (10% condition,

**Shape** of the sampling distribution of
       is approximately normal when...
1) Both population distributions are Normal
2) When both sample sizes are large enough
   ($n_1 \geq 30$ and $n_2 \geq 30$)

3) Small samples —
   You must graph (histograms) BOTH samples
   to look for skewness and/or severe outliers

# ■ The Two-Sample $t$ Statistic

## Important Formula:

When data come from two random samples or two groups in a randomized experiment, the statistic $\bar{x}_1 - \bar{x}_2$ is our best guess for the value of $\mu_1 - \mu_2$.

When the Independent condition is met, the standard deviation of the statistic $\bar{x}_1 - \bar{x}_2$ is :

$$\sigma_{\bar{x}_1 - \bar{x}_2} = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

Since we don't know the values of the parameters $\sigma_1$ and $\sigma_2$, we replace them in the standard deviation formula with the sample standard deviations. The result is the **standard error** of the statistic $\bar{x}_1 - \bar{x}_2$ :

$$\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

If the Normal condition is met, **we standardize the observed difference to obtain a $t$ statistic** that ...served difference is from its mean in standard devia...

$$t = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

**The two-sample $t$ statistic has approximately a $t$ distribution.** We can use technology to determine degrees of freedom OR we can use a conservative approach, using the smaller of $n_1 - 1$ and $n_2 - 1$ for the degrees of freedom.

DF

*TEST STATISTIC FOR 2 SAMPLE MEANS is "$t$".*

*2 WAYS TO DETERMINE "DF":*
*① TECHNOLOGY - Calculator DF*
*② CONSERVATIVE - DF is based on the smaller sample size.*
*✱ MUST STATE METHOD USED!!*

## ■ Confidence Intervals for $\mu_1 - \mu_2$

### Two-Sample $t$ Interval for a Difference Between Means

When the Random, Normal, and Independent conditions are met, an approximate level C confidence interval for $(\bar{x}_1 - \bar{x}_2)$ is

$$\left[(\bar{x}_1 - \bar{x}_2) \pm t^* \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}\right.$$

*Formula FOR 2 sample CI.*

where $t^*$ is the critical value for confidence level C for the $t$ distribution with degrees of freedom from either technology or the smaller of $n_1 - 1$ and $n_2 - 1$.

## ■ Significance Tests for $\mu_1 - \mu_2$

### Two-Sample $t$ Test for the Difference Between Two Means

If the Random, Normal, and Independent conditions are met, we can proceed:
To do a test, standardize $\bar{x}_1 - \bar{x}_2$ to get a two-sample $t$ statistic:

$$\text{test statistic} = \frac{\text{statistic} - \text{parameter}}{\text{standard deviation of statistic}}$$

$$t = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

To find the $P$-value, use the $t$ distribution with degrees of freedom given by technology or by the conservative approach (df = smaller of $n_1 - 1$ and $n_2 - 1$).

*2 Sample $t$ statistic for T.O.H.*

**EXAMPLE:** Based on information from the U.S. National Health and Nutrition Examination Survey (NHANES), the heights (in inches) of ten-year-old girls follow a Normal distribution $N(56.4, 2.7)$. The heights (in inches) of ten-year-old boys follow a Normal distribution $N(55.7, 3.8)$. A researcher takes independent SRSs of 12 girls and 8 boys of this age and measures their heights. After analyzing the data, the researcher reports that the sample mean height of the boys is larger than the sample mean height of the girls

a) Describe the shape, center, and spread of the sampling distribution of $\bar{x}_f - \bar{x}_m$.

b) Find the probability of getting a difference in sample means $\bar{x}_f - \bar{x}_m$ that is less than 0.

c) Does the result in part (b) give us reason to doubt the researchers' stated results?

---

## 2 POPULATIONS

FEMALE HEIGHTS $\quad \mu_F = 56.4 \text{ in} \quad N(56.4, 2.7)$

MALE HEIGHTS $\quad \mu_m = 55.7 \quad\quad N(55.7, 3.8)$

SRS
$n_1 = 12$
$n_2 = 8$

---

(A) THE SAMPLING DISTRIBUTION OF $\bar{X}_F - \bar{X}_m$

① SHAPE: Since both population distributions are normal, THE SAMPLING DISTRIBUTION OF $\bar{X}_F - \bar{X}_m$ IS APPROXIMATELY NORMAL

② CENTER:

$$\underbrace{\mu_{\bar{X}_F - \bar{X}_m}}_{\text{Sampling Distribution}} = \underbrace{\mu_F - \mu_m}_{\text{Pop. Parameters}} = 56.4 - 55.7$$

$$= \boxed{0.7 \text{ INCHES}}$$

③ SPREAD:

$$\sigma_{\bar{X}_F - \bar{X}_m} = \sqrt{\frac{\sigma_F^2}{n_F} + \frac{\sigma_M^2}{n_M}}$$

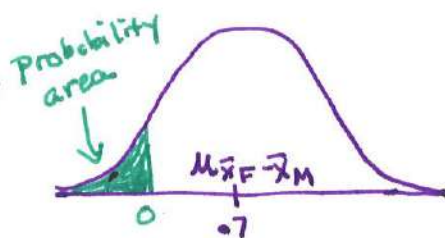$$= \sqrt{\frac{2.7^2}{12} + \frac{3.8^2}{8}} = \boxed{1.55 \text{ inches}}$$

**B** FIND THE PROBABILITY OF GETTING
 DIFFERENCE $(\bar{X}_F - \bar{X}_M)$ LESS THAN 0.

WRITE A
PROBABILITY
STATEMENT →

$$P(\bar{X}_F - \bar{X}_M < 0) = \underline{.3258}$$

THIS IS EQUIVALENT TO $P(\bar{X}_F < \bar{X}_M)$
THAT IS "MEAN HT OF 10 YEAR OLD
GIRL IS SHORTER THAN A BOY"
OR "mean HT OF A BOY IS
TALLER THAN A GIRL" $P(\bar{X}_M > \bar{X}_F)$

Graph the
Sampling
distribution →

Probability
area



$\mu_{\bar{X}_F - \bar{X}_M}$
0   .7

TO FIND THE PROBABILITY:
normalcdf$(-E99, 0, .7, 1.55) =$
LB   UB   $\mu$   $\sigma$   $\boxed{.3258}$

OR FIND PROBABILITY BY STANDARDIZING $Z$:

$N(0,1)$ →  $Z = \dfrac{0 - .7}{1.55} = -.45$

$P(Z \le -.45) = \boxed{.3264}$    normal cdf $(-E99, -.45, 0, 1)$

**C** Researcher claims boys are taller than
girls (at 10 years old)

Based on these results, there is about a 33%
Chance of getting a sample mean difference
less than zero (0) due to sampling variability.

THIS MEANS THAT WE WOULD EXPECT
1 OUT 3 10 year old boys to be taller than girls.
Since this is not an unusual result, we
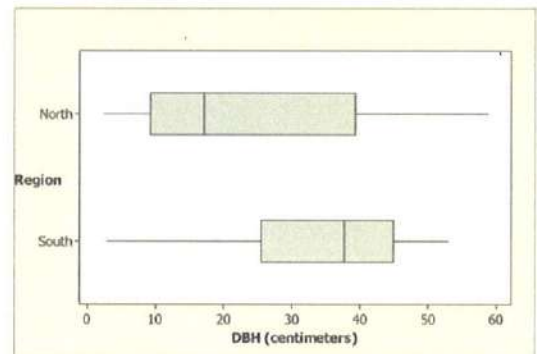Should not doubt the researchers claim that
boys are taller than girls.

■ **Example:** Big Trees, Small Trees, Short Trees, Tall Trees

The Wade Tract Preserve in Georgia is an old-growth forest of longleaf pines that has survived in a relatively undisturbed state for hundreds of years. One question of interest to foresters who study the area is "How do the sizes of longleaf pine trees in the northern and southern halves of the forest compare?" To find out, researchers took random samples of 30 trees from each half and measured the diameter at breast height (DBH) in centimeters. Comparative boxplots of the data and summary statistics from Minitab are shown below. Construct and interpret a 90% confidence interval for the difference in the mean DBH for longleaf pines in the northern and southern halves of the Wade Tract Preserve.

Descriptive Statistics: North, South

| Variable | N | Mean | StDev |
|----------|-----|--------|-------|
| $\mu_2$ North | 30 | 23.70 ✔ | 17.50 |
| $\mu_1$ South | 30 | 34.53 ✔ | 14.26 |



**NAME OF INTERVAL:**

2 Sample t-interval for the difference of means $(\mu_1 - \mu_2)$

TIP: look at the sample means and make $\mu_1$ the larger mean to have a positive difference.

**DEFINE PARAMETERS:**

$\mu_1 =$ TRUE mean diameter of trees in the SOUTH    $\bar{x}_1 = 34.53$
$\mu_2 =$ TRUE mean diameter of trees in the NORTH    $\bar{x}_2 = 23.70$

**SIGNIFICANCE LEVEL:**

Want to estimate the difference $(\mu_1 - \mu_2)$ at the 90% C.L.    $\alpha = .10$

**CONDITIONS:**

$\sigma$ IS UNKNOWN (t-interval)

RANDOM — Random samples from both the NORTH and SOUTH

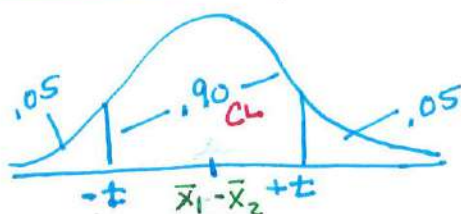INDEPENDENT — ① The samples from the North and South are independent.
② IT IS FAIR TO ASSUME THAT THERE ARE AT LEAST 10(30) = 300 TREES IN EACH REGION.

NORMAL — SINCE BOTH SAMPLE SIZES ARE RELATIVELY LARGE WITH SAMPLE SIZES OF 30, IT IS REASONABLE BOTH DISTRIBUTIONS ARE APPROXIMATELY NORMAL.

# CONSTRUCT A 2 SAMPLE T INTERVAL FOR $\mu_1 - \mu_2$

$$\boxed{\alpha = .10}$$

.05    .90 CL    .05

$-t \quad \bar{x}_1 - \bar{x}_2 \quad +t$

## HAND CALCULATE

FORMULA:

$$(\bar{x}_1 - \bar{x}_2) \pm t^* \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

Conservative $df = n-1$ of smaller of the 2 samples

$$\boxed{df = 29}$$

$t^* = \text{invT}(.05, 29)$    $t^* = \pm 1.70$
     ↑ DF

$34.53 - 23.70 \pm (1.70) \sqrt{\frac{14.26^2}{30} + \frac{17.50^2}{30}}$

$10.83 \pm (1.70)(4.12)$   ← SE

$10.83 \pm 7.01$

    ↑ ME    $(3.82, 17.84)$

## CALCULATOR

### 2-Samp T INTERVAL

| SOUTH | NORTH |
|---|---|
| $\bar{x}_1 = 34.53$ | $\bar{x}_2 = 23.70$ |
| $s_{x_1} = 14.26$ | $s_{x_2} = 17.50$ |
| $n_1 = 30$ | $n_2 = 30$ |

POOLED NO ⟵ ALWAYS
=
↓

$(3.9362, 17.724)$

$\boxed{df = 55.7}$

TECHNOLOGY. COMPLEX FORMULA
DO NOT NEED TO KNOW

\* YOU **must** state what df
you used (Conservative OR
TECHNOLOGY)!!!

## CONCLUDE:

WE ARE 90% CONFIDENT THAT THE INTERVAL 3.83 to 17.83 cm CAPTURES THE TRUE DIFFERENCE IN THE ACTUAL MEAN DBH BETWEEN THE SOUTHERN AND NORTHERN TREES.

THIS SUGGESTS THAT THE MEAN DIAMETER OF SOUTHERN TREES IS BETWEEN 3.83 AND 17.83 cm LARGER THAN THE MEAN DIAMETER OF THE NORTHERN TREES, ON AVERAGE.
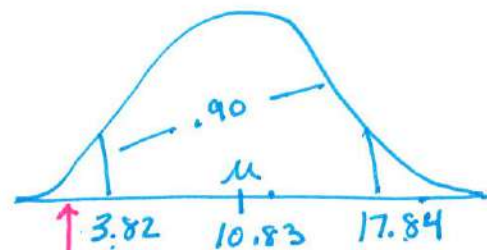
## DISCUSSION

BASED ON THE CI, IS THERE CONVINCING EVIDENCE THAT THE TREE DIAMETER IS DIFFERENT BETWEEN THE NORTH AND SOUTH?

$H_0: \mu_N = \mu_S$

$H_A: \mu_N \neq \mu_S$

.90

$\mu$

3.82    10.83    17.84

$\alpha = .10$

CI

Reject    FAIL TO    Reject
$H_0$    REJECT    $H_0$

SINCE 0 IS NOT IN THE CI, We have convincing evidence to reject $H_0$ and believe the diameters in North & South are different.

SIGNIFICANCE LEVEL
$\alpha = .10$   10% Chance of making a TYPE 1 ERROR

9

■ **Example:** Calcium and Blood Pressure

Does increasing the amount of calcium in our diet reduce blood pressure? Examination of a large sample of people revealed a relationship between calcium intake and blood pressure. The relationship was strongest for black men. Such observational studies do not establish causation. Researchers therefore designed a randomized comparative experiment. The subjects were 21 healthy black men who volunteered to take part in the experiment. They were randomly assigned to two groups: 10 of the men received a calcium supplement for 12 weeks, while the control group of 11 men received a placebo pill that looked identical. The experiment was double-blind. The response variable is the decrease in systolic (top number) blood pressure for a subject after 12 weeks, in millimeters of mercury. An increase appears as a negative response Here are the data:

| Group 1 (calcium): | 7 | −4 | 18 | 17 | −3 | −5 | 1 | 10 | 11 | −2 | |
| Group 2 (placebo): | −1 | 12 | −1 | −3 | 3 | −5 | 5 | 2 | −11 | −1 | −3 |

CALCIUM + BLOOD PRESSURE

## 2 SAMPLE $t$-test for $\mu_1 - \mu_2$

GRAPHS

**GROUP 1 (calcium)** $\quad \overline{X}_1 = 5.0 \qquad S_1 = 8.74 \qquad n_1 = 10$

**GROUP 2 (placebo)** $\quad \overline{X}_2 = -.27 \quad S_2 = 5.90 \qquad n_2 = 11$

**PARAMETERS:** $\mu_1 =$ true mean decrease in blood pressure (Calcium Supplements)

$\mu_2 =$ true mean decrease in blood pressure (PLACEBO)

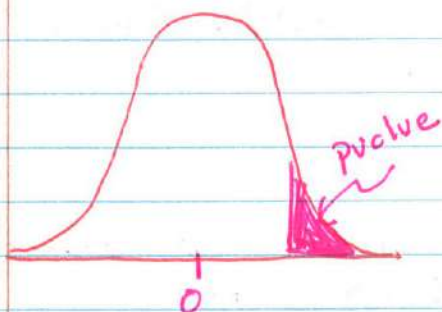**HYPOTHESIS**

$H_0 : \mu_1 - \mu_2 = 0$
$H_A : \mu_1 - \mu_2 > 0$

$H_0 : \mu_1 = \mu_2$
$H_A : \mu_1 > \mu_2$

SIGNIFICANCE LEVEL: $\alpha = .05$

**SKETCH GRAPH:**



pvalue

0

Placebo



Calcium



-15     0     15

blood pressure

**CONDITIONS:** Random, Normal Independent, $\sigma$
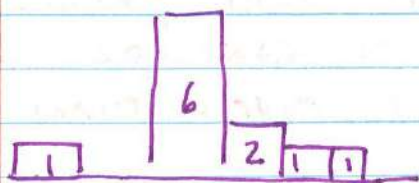
• Random — 21 subjects were randomly assigned TO THE 2 TREATMENTS.

INDEPENDENT: ① Due to Random assignment, these 2 groups can be viewed as independent.
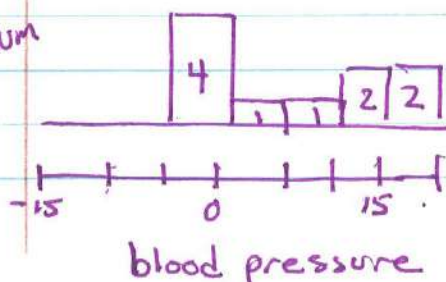② Reasonable individual observations are independent

Normal: Since samples are both under 30, We looked at graphs (above) and do not show clear evidence of skewness & no outliers

$\sigma$ UNKNOWN ($t$ inference)

# STATE TEST BY NAME OR FORMULA

NAME: **2 Sample T test for $\mu_1 - \mu_2$**

TEST STATISTIC $\quad t = \dfrac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{s_1^2/n_1 + s_2^2/n_2}}$

STATE DF: DF = 9 (conservative)

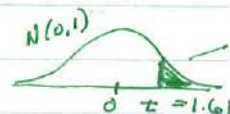$$t = \frac{(5 - (-.27)) - 0}{SE \rightsquigarrow \sqrt{\frac{8.74^2}{10} + \frac{5.90^2}{11}}} = \frac{5.27}{3.28} = 1.61 \qquad \boxed{t = 1.61}$$

Pvalue =

**PVALUE**   State probability: $P(t \geq 1.61) = .071$
find pvalue
write decision in conclusion



$\curvearrowleft$ tcdf $(1.61, E99, 9)$

Check w/ CALC: (STAT) (TEST) 2 SAMP T TEST
 * ALWAYS USE Pooled $\boxed{NO}$ $\leftarrow$ we are not going to pool variances

What you need to write using calc.
$\left[\begin{array}{l} t = \underline{1.60} \quad df = \underline{15.6} \quad p\text{value} = P(t \geq 1.60) = .064 \\ \qquad\qquad \text{(technology)} \end{array}\right.$

**CONCLUSION** BECAUSE THE PVALUE IS GREATER THAN $\alpha = .05$, WE FAIL TO REJECT Ho

THE EXPERIMENT DID NOT PROVIDE CONVINCING EVIDENCE TO CONCLUDE CALCIUM REDUCES BLOOD PRESSURE MORE THAN A PLACEBO

# STATE TEST BY NAME OR FORMULA

NAME: 2 Sample T test for $\mu_1 - \mu_2$

TEST STATISTIC $t = \dfrac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{s_1^2/n_1 + s_2^2/n_2}}$

STATE DF:

---

## PVALUE

State probability:
find pvalue
write decision in conclusion

Check w/ CALC: (STAT) (TEST) 2 SAMP T TEST
  * ALWAYS USE Pooled NO ← we are not going
                                      to pool variances

$t =$ ___ $df =$ ___ $p =$ ___

---

## CONCLUSION

BECAUSE THE PVALUE IS GREATER THAN $\alpha = .05$, WE FAIL TO REJECT Ho

THE EXPERIMENT DID NOT PROVIDE CONVINCING EVIDENCE TO CONCLUDE CALCIUM REDUCES BLOOD PRESSURE MORE THAN A PLACEBO

| EXAMPLE | CALCIUM + BLOOD PRESSURE |

$$2 \text{ SAMPLE } t\text{-test for } \mu_1 - \mu_2$$

GRAPHS

GROUP 1
(calcium)

$\overline{X}_1 =$    $S_1 =$    $n_1 =$

GROUP 2
(placebo)

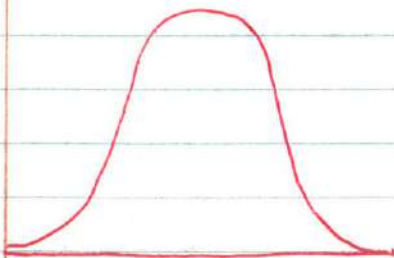$\overline{X}_2 =$    $S_2 =$    $n_2 =$

PARAMETERS:    $\mu_1 =$

$\mu_2 =$

HYPOTHESIS

SIGNIFICANCE
LEVEL:

SKETCH
GRAPH:



CONDITIONS: Random, Normal
Independent, $\sigma$

• Random − 21 subjects were
randomly assigned TO THE
2 TREATMENTS.

INDEPENDENT: ① Due to Random
assignment, these 2 groups can be
viewed as independent.
② Reasonable individual observations
are independent

Normal: Since samples are both
under 30, We looked at graphs
(above) and do not show clear
evidence of skewness & No outliers

$\sigma$ UNKNOWN (t inference)