

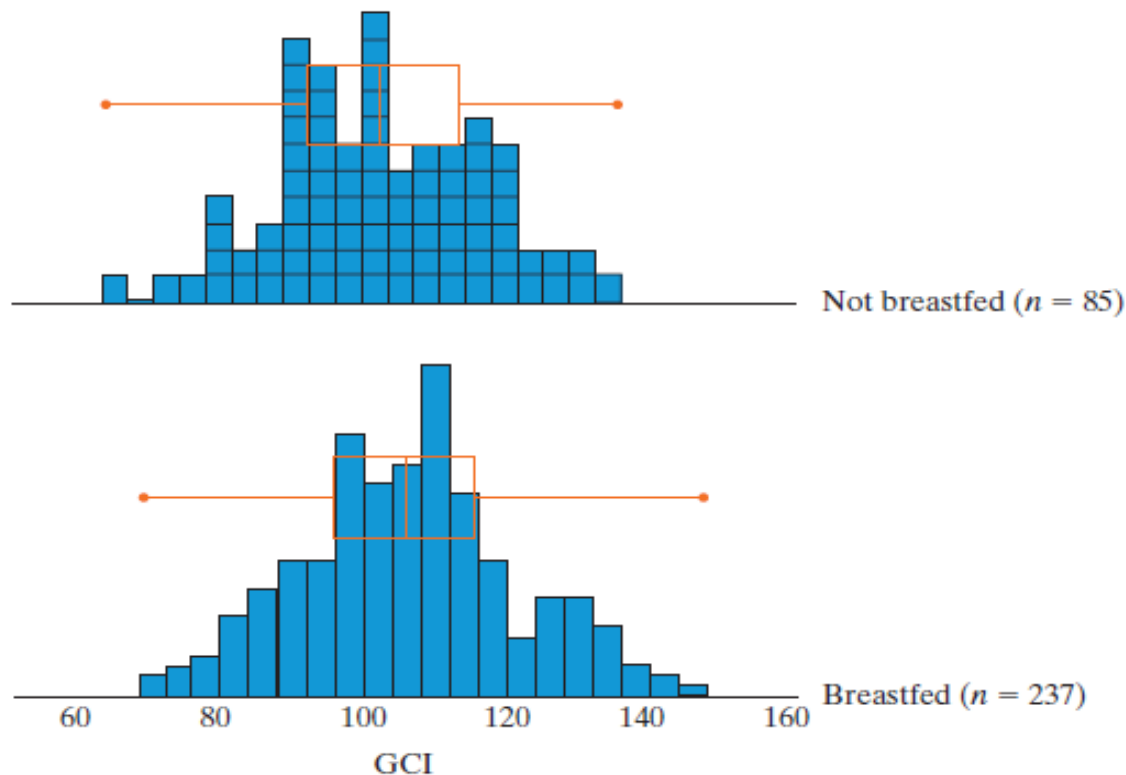
## Stat 13, Intro. to Statistical Methods for the Life and Health Sciences.

1. Breastfeeding and intelligence example continued.
2. Paired data, studying with music, running bases,
3. When to use which formula.
4. Multiple testing and publication bias.

Read through ch7.

# Breastfeeding and Intelligence

Group	Sample size, $n$	Sample mean	Sample SD
Breastfed	237	105.3	14.5
Not BF	85	100.9	14.0



# T-statistic

- If we can assume the draws are iid and the populations are normal, with unknown sds, then t-statistic is used.
- It is the number of standard deviations our statistic is above or below the mean under the null hypothesis.

- $$t = \frac{\text{statistic} - \text{hypothesized value}}{SE} = \frac{\bar{x}_1 - \bar{x}_2 - 0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

- Here, 
$$t = \frac{105.3 - 100.9}{\sqrt{\left(\frac{14.5^2}{237} + \frac{14.0^2}{85}\right)}} = 2.46. \text{ p-value} \sim 1.4\%.$$

- $2 * \text{pnorm}(2.46, \text{lower} = \text{F})] = 1.39\%,$   
 $2 * \text{pt}(2.46, \text{lower} = \text{F}, \text{df} = 320) = 1.44\%. \text{ df} = n_1 + n_2 - 2 \text{ here.}$

# Breastfeeding and Intelligence

Meaning of the p-value:

- If breastfeeding were not related to GCI at age 4, then the probability of observing a difference of 4.4 or more or -4.4 or less just by chance is about 1.4%.

- A 95% CI can also be obtained using the t-distribution. The SE is  $\sqrt{\left(\frac{14.5^2}{237} + \frac{14.0^2}{85}\right)} = 1.79$ .  
So the margin of error is multiplier x SE.

# Breastfeeding and Intelligence

- The SE is  $\sqrt{\left(\frac{14.5^2}{237} + \frac{14.0^2}{85}\right)} = 1.79$ . The margin of error is multiplier x SE.
- The multiplier should technically be obtained using the t distribution, but for large sample sizes you get almost the same multiplier with t and normal. Use 1.96 for a 95% CI to get  $4.40 \pm 1.96 \times 1.79 = 4.40 \pm 3.51 = (0.89, 7.91)$ .
- The book uses 2 instead of 1.96, and the applet uses 1.9756 from the t-distribution. Just use 1.96 for this class.

# Breastfeeding and Intelligence

- We have strong evidence against the null hypothesis and can conclude the association between breastfeeding and intelligence here is statistically significant.
- Breastfed babies have statistically significantly higher average GCI scores at age 4.
- We can see this in both the small p-value (0.015) and the confidence interval that says the mean GCI for breastfed babies is 0.89 to 7.91 points higher than that for non-breastfed babies.

# Breastfeeding and Intelligence

- Can you conclude that breastfeeding improves average GCI at age 4?
  - No. The study was not a randomized experiment.
  - We cannot conclude a cause-and-effect relationship.
- There might be alternative explanations for the significant difference in average GCI values.
- What might some confounding factors be?

# Breastfeeding and Intelligence

- Can you conclude that breastfeeding improves average GCI at age 4?
  - No. The study was not a randomized experiment.
  - We cannot conclude a cause-and-effect relationship.
- There might be alternative explanations for the significant difference in average GCI values.
  - Maybe better educated mothers are more likely to breastfeed their children
  - Maybe mothers that breastfeed spend more time with their children and interact with them more.
  - Some mothers who do not breastfeed are less healthy or their babies have weaker appetites and this might slow down development in general.



# Breastfeeding and Intelligence

- Could you design a study that allows drawing a cause-and-effect conclusion?
  - We would have to run an experiment using random assignment to determine which mothers breastfeed and which would not. (It would be impossible to double-blind.)
  - Random assignment roughly balances out all other variables.
- Is it feasible/ethical to conduct such a study?

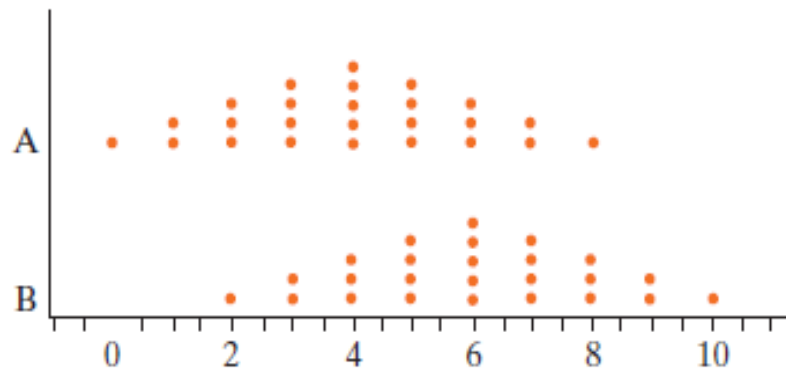
# Strength of Evidence

- We already know:
  - As sample size increases, the strength of evidence increases.
  - Just as with proportions, as the effect size increases, the strength of evidence increases.

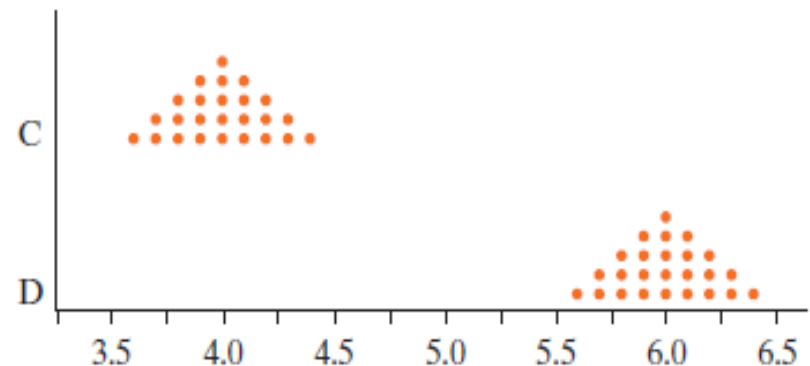
# More Strength of Evidence

- If the effect size is held constant, i.e. the means are the same distance apart, but the standard deviations decrease, then the strength of evidence increases.
- Which gives stronger evidence against the null?

Dotplot Pair 1



Dotplot Pair 2



- Smaller SDs lead to stronger evidence against the null.

# Effects on Width of Confidence Intervals

- Just as before:
  - As sample size increases, confidence interval widths tend to decrease.
  - As confidence level increases, confidence interval widths increase.
  - The effect size, i.e. the difference in means, will not affect the width (margin of error) but will affect the center of the CI.
- As we saw with a single mean, as the SDs of the samples increase, the width of the confidence interval will increase.

# Paired Data.

Chapter 7

# Introduction

- The paired data sets in chapter 7 have one *pair* of quantitative response values for each obs. unit.
- This allows for a comparison where the other possible confounders are as similar as possible between the two groups.
  - The big idea is with paired data, just view the *differences* between each pair of scores as your data. Now you have one variable and can just analyze it using the methods we already know.
  - When you analyze paired data this way, person to person variability gets removed so you get more power when testing, smaller p-values and smaller margins of error.

# Paring and Observational Studies

**You can often do matched pairs in observational studies, when you know the potential confounder ahead of time.**

If you are studying whether the portacaval shunt decreases the risk of heart attack, you could match each patient getting the shunt with a patient of similar health not getting the shunt.

If you are studying whether lefthandedness causes death, and you want to account for age in the population, you could match each leftie with a rightie of the same age, and compare their ages at death.

# Simulation-Based Approach for Analyzing Paired Data, and rounding first base example.

Section 7.2

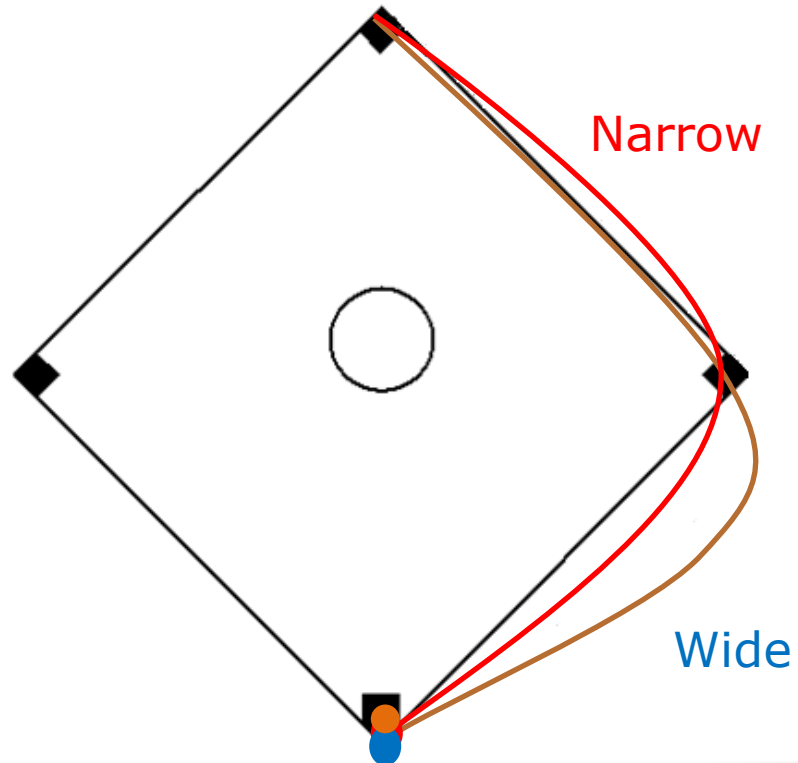


# Rounding First Base

Example 7.2

# Rounding First Base

- Imagine you've hit a line drive and are trying to reach second base.
- Does the path that you take to round first base make much of a difference?
  - **Narrow angle**
  - **Wide angle**



# Rounding First Base

- Woodward (1970) investigated these base running strategies.
- He timed 22 different runners from a spot 35 feet past home to a spot 15 feet before second.
- Each runner used each strategy (paired design), with a rest in between.
- He used random assignment to decide which path each runner should do first.
- **This paired design controls for the runner-to-runner variability.**

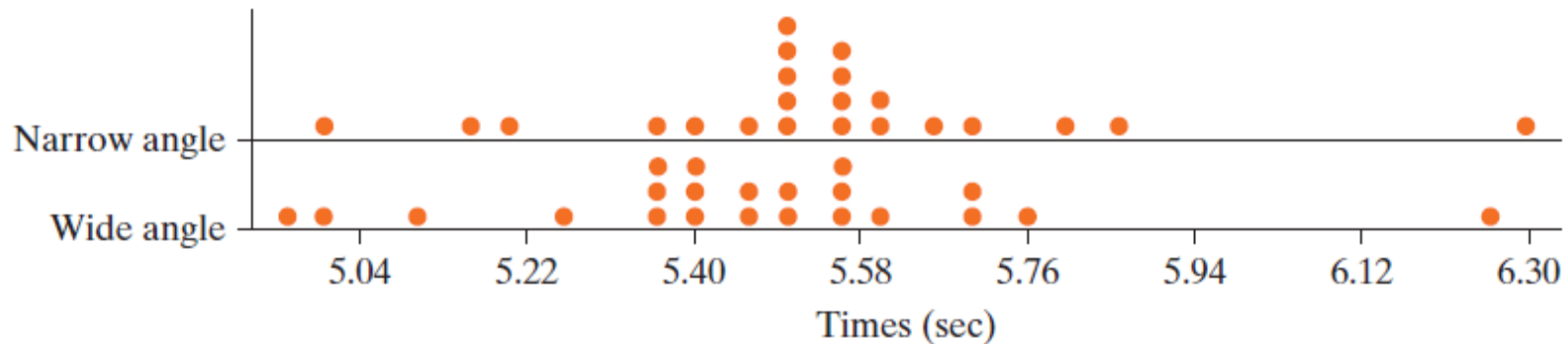
# First Base

- What are the observational units in this study?
  - The runners (22 total)
- What variables are recorded? What are their types and roles?
  - Explanatory variable: base running method: wide or narrow angle (categorical)
  - Response variable: time from home plate to second base (quantitative)
- Is this an observational study or an experiment?
  - Randomized experiment.

# The results

**TABLE 7.1** The running times (seconds) for the first 10 of the 22 subjects

Subject	1	2	3	4	5	6	7	8	9	10	
Narrow angle	5.50	5.70	5.60	5.50	5.85	5.55	5.40	5.50	5.15	5.80	...
Wide angle	5.55	5.75	5.50	5.40	5.70	5.60	5.35	5.35	5.00	5.70	...



# Paired data and rounding first base example.

- There is a lot of overlap in the distributions and substantial variability.

	Mean	SD
Narrow	5.534	0.260
Wide	5.459	0.273

- It is difficult to detect a difference between the methods when there is so much variation.
-

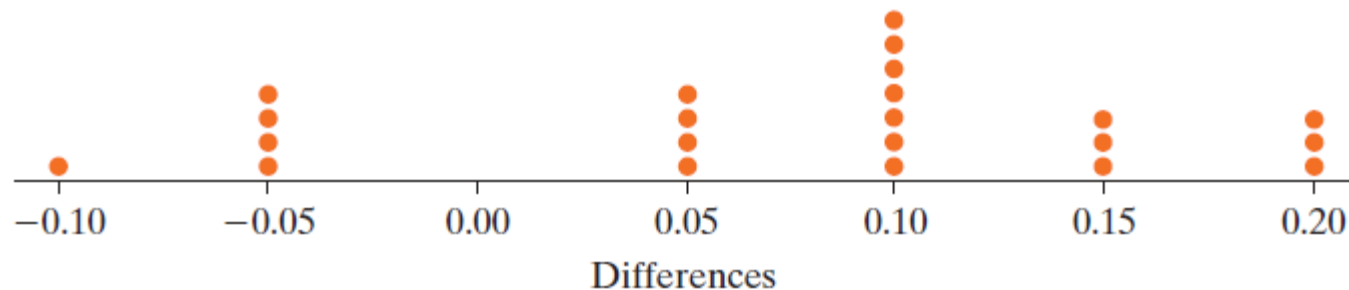
# Rounding First Base

- These data are clearly paired.
- The paired response variable is time difference in running between the two methods and we can use this in analyzing the data.

# The Differences in Times

**TABLE 7.2** Last row is difference in times for each of the first 10 runners (narrow – wide)

Subject	1	2	3	4	5	6	7	8	9	10	
Narrow angle	5.50	5.70	5.60	5.50	5.85	5.55	5.40	5.50	5.15	5.80	...
Wide angle	5.55	5.75	5.50	5.40	5.70	5.60	5.35	5.35	5.00	5.70	...
Difference	-0.05	-0.05	0.10	0.10	0.15	-0.05	0.05	0.15	0.15	0.10	...



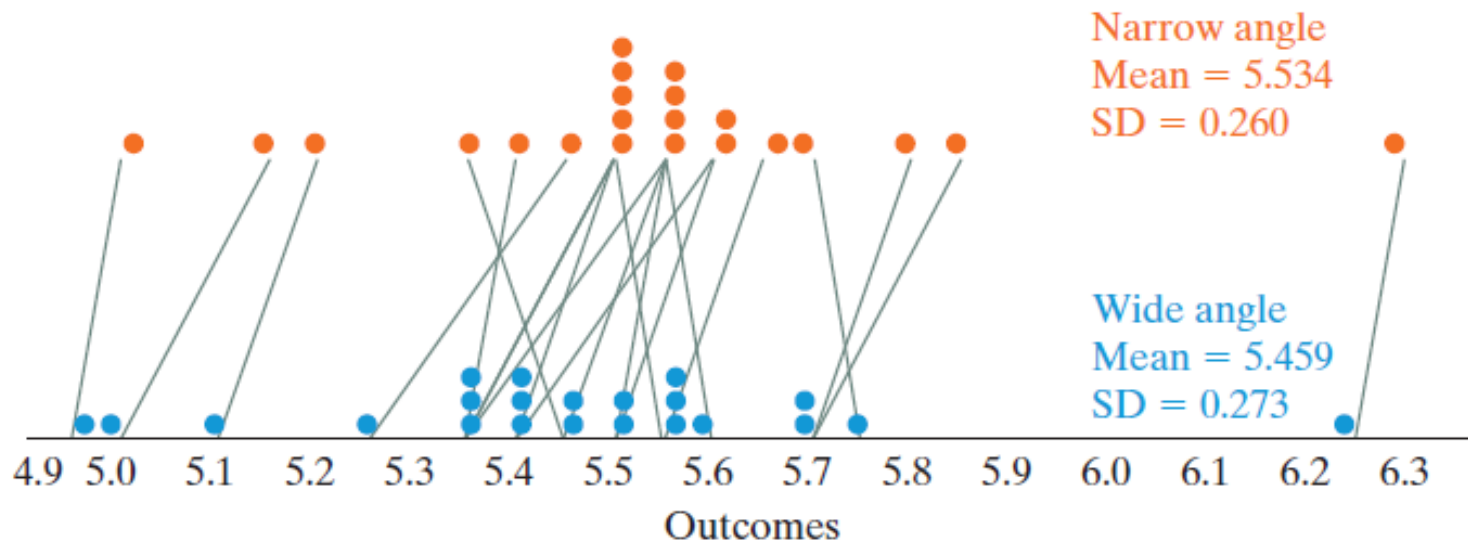


# The Differences in Times

- Mean difference is  $\bar{x}_d = 0.075$  seconds
- Standard deviation of the differences is  $SD_d = 0.0883$  sec.
- This standard deviation of 0.0883 is smaller than the original standard deviations of the running times, which were 0.260 and 0.273.

# Rounding First Base

- Below are the original dotplots with each observation paired between the base running strategies.
- What do you notice?



# Rounding First Base

- Is the average difference of  $\bar{x}_d = 0.075$  seconds significantly different from 0?
- The parameter of interest,  $\mu_d$ , is the long run mean difference in running times for runners using the narrow angled path instead of the wide angled path. (narrow – wide)

# Rounding First Base

The hypotheses:

- $H_0: \mu_d = 0$ 
  - The long run mean difference in running times is 0.
- $H_a: \mu_d \neq 0$ 
  - The long run mean difference in running times is not 0.
- The statistic  $\bar{x}_d = 0.075$  is above zero.
- How likely is it to see an average difference in running times this big or bigger by chance alone, even if the base running strategy has no genuine effect on the times?

# Rounding First Base

How can we use simulation-based methods to find an approximate p-value?

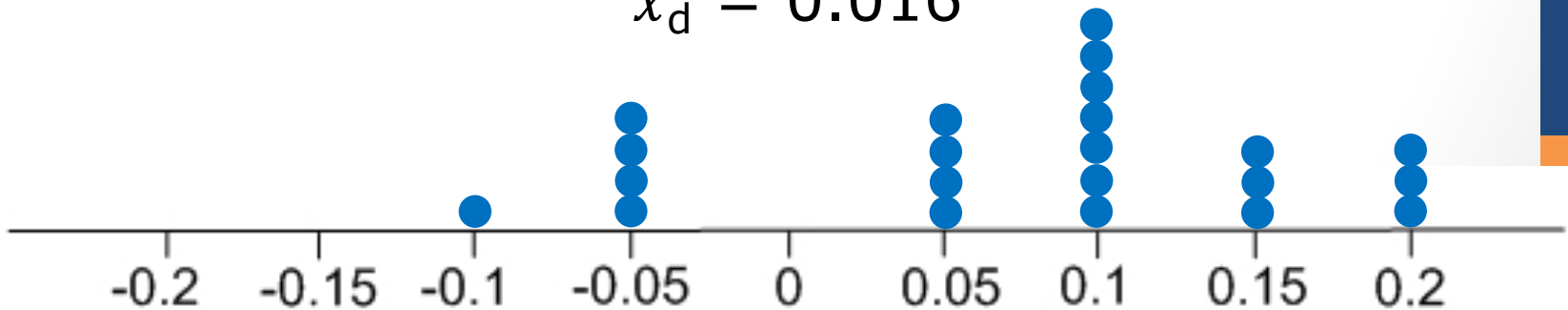
- The null hypothesis says the running path does not matter.
- So we can use our same data set and, for each runner, randomly decide which time goes with the narrow path and which time goes with the wide path and then compute the difference. (Notice we do not break our pairs.)
- After we do this for each runner, we then compute a mean difference.
- We will then repeat this process many times to develop a null distribution.

# Random Swapping

Subject	1	2	3	4	5	6	7	8	9	10	
narrow angle	5.50	5.70	5.60	5.50	5.85	5.55	5.40	5.50	5.15	5.80	...
wide angle	5.55	5.75	5.50	5.40	5.70	5.60	5.35	5.35	5.00	5.70	...
diff	0.05	-0.05	-0.10	0.10	0.15	0.05	0.05	0.15	0.15	-0.10	...

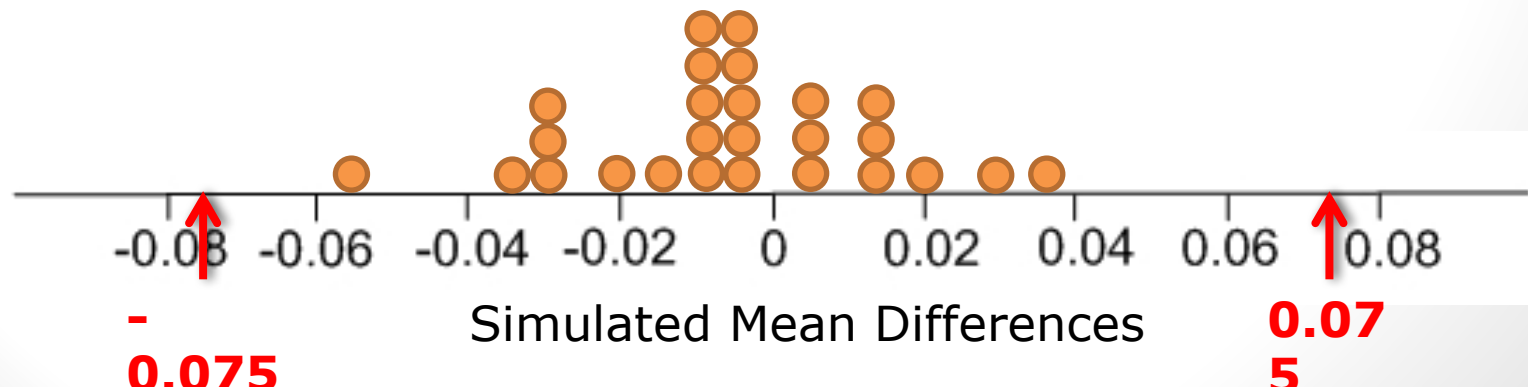


$$\bar{x}_d = 0.016$$



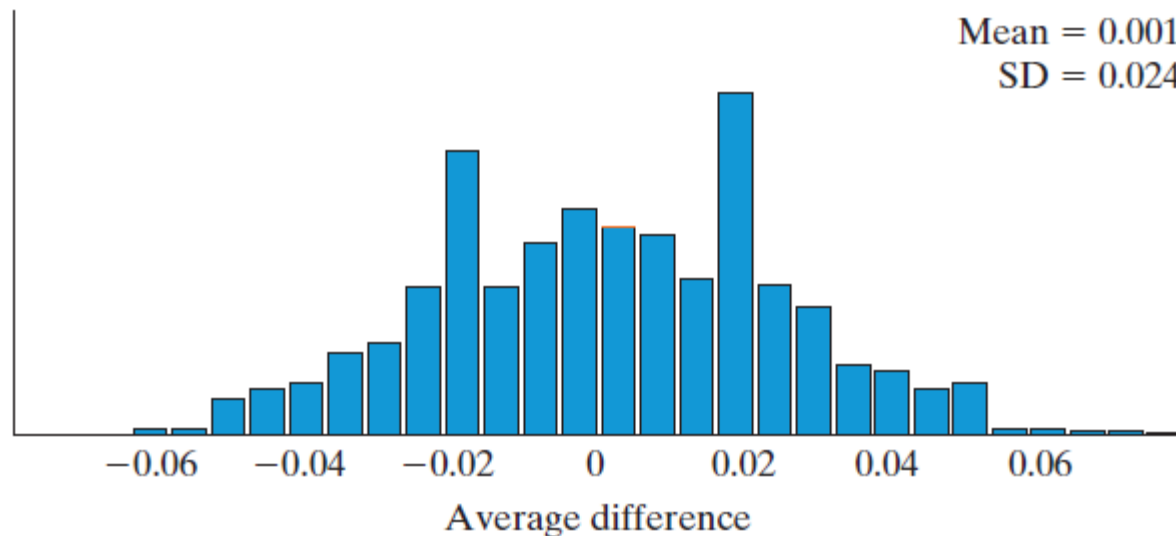
# More Simulations

With 26 repetitions of creating simulated mean differences, we did not get any that were as extreme as 0.075.



# First Base

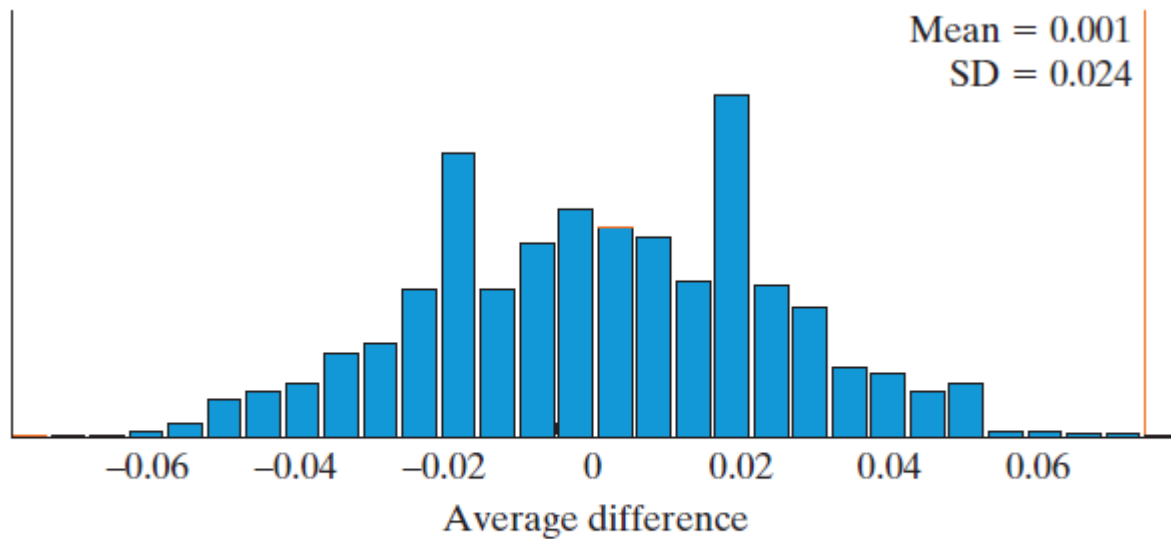
- Here is a null distribution of 1000 simulated mean differences.
- Notice it is centered at zero, which makes sense in agreement with the null hypothesis.
- Notice also the SD of these MEAN DIFFERENCES is  $0.024 = SE$ . SD of time differences was 0.0883. SD of mean time diff.s = .024.
- Where is our observed statistic of 0.075?





# First Base

- Only 1 of the 1000 repetitions of random swappings gave a  $\bar{x}_d$  value at least as extreme as 0.075.

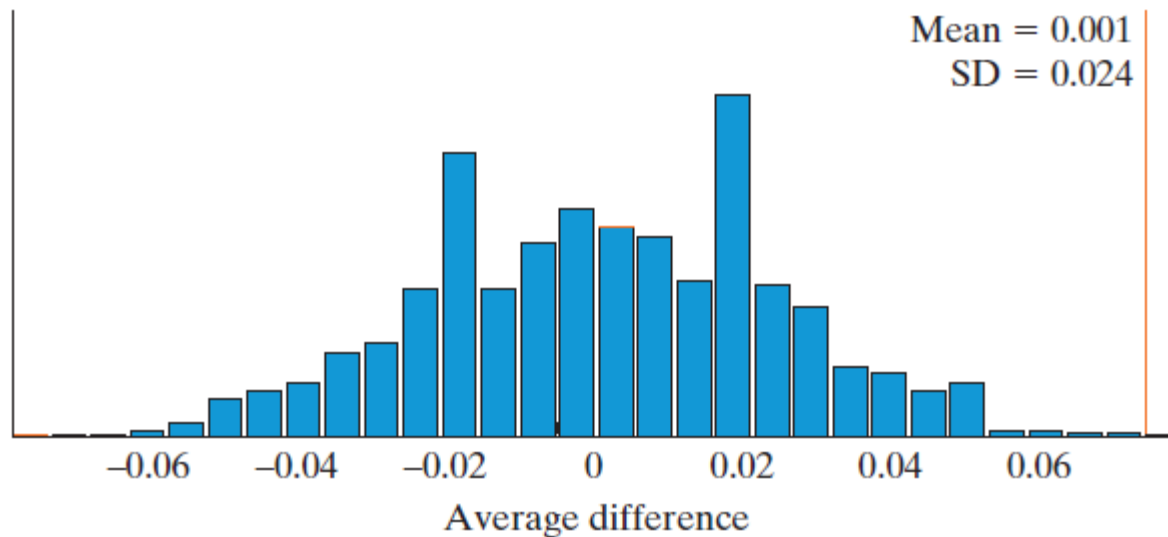


Count samples:

Count = 1/1000 (0.0010)

# First Base

- We can also standardize 0.075 by dividing by the SE of 0.024 to see our standardized statistic =  $\frac{0.075}{0.024} = 3.125$ .



Count samples:

Count = 1/1000 (0.0010)

# First Base

- The simulation p-value is 0.1%. We can also standardize 0.075 by dividing by the SE of 0.024 to see our standardized statistic =  $\frac{0.075}{0.024} = 3.125$ .
- If we had used the formula instead of simulations to get the SE, we would have obtained  $s/\sqrt{n} = .0883/\sqrt{22} = .019$  instead of .024. Here  $s = .0883$ .  $s$  is the sample sd of the differences. Using the formula for the SE, we would get a standardized statistic of  $.075/.019 = 3.95$ .
- Either way, clearly stat. sig. Using the formulas, you would get  $2 * pt(3.95, lower=F, df=21) = 0.0732\%$  instead of 0.1%. This relies on assuming the time differences are normal though.

# Rounding First Base

- With a p-value of 0.1%, we have very strong evidence against the null hypothesis. The running path makes a statistically significant difference with the wide-angle path being faster on average.
- We can draw a cause-and-effect conclusion since the researcher used random assignment of the two base running methods for each runner.
- There was not much information about how these 22 runners were selected though so it is unclear if we can generalize to a larger population.

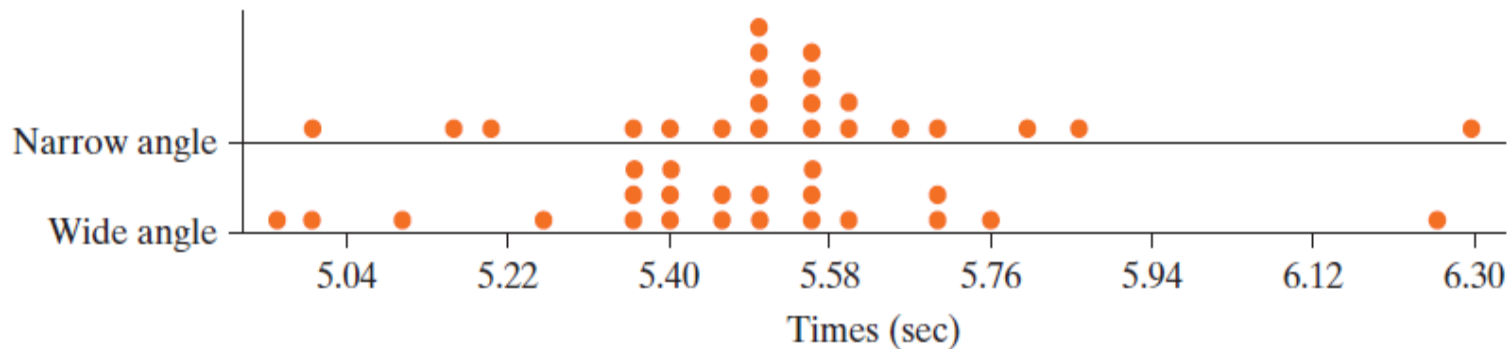
# First Base

- Using the simulation based SE of .024, approximate a 95% confidence interval for  $\mu_d$ :
  - $0.075 \pm 1.96(0.024)$  seconds.
  - $(0.028, 0.122)$  seconds.
- What does this mean?
  - We are 95% confident that, if we were to keep testing this indefinitely, the narrow angle route would take somewhere between 0.028 to 0.122 seconds longer on average than the wide angle route.

# First Base

## Alternative Analysis

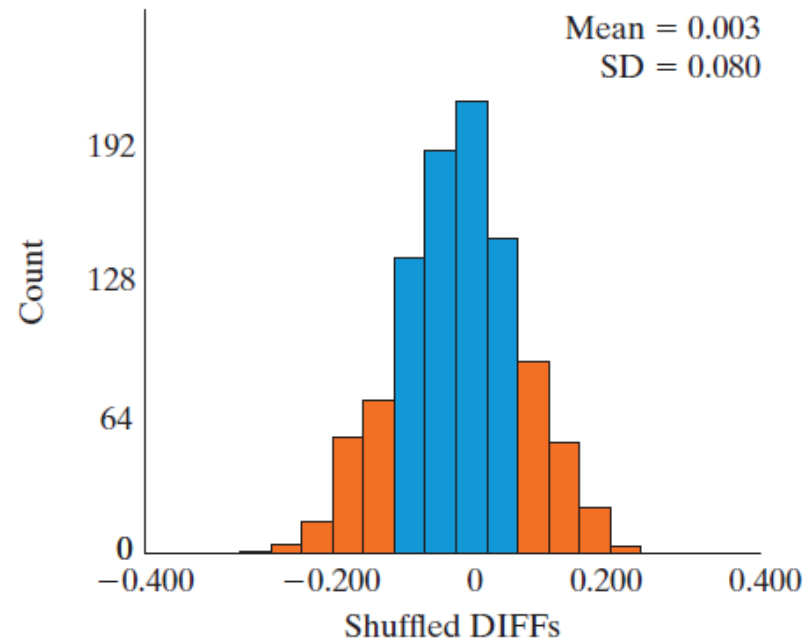
- What do you think would happen if we wrongly analyzed the data using a 2 independent samples procedure? (i.e. the researcher selected 22 runners to use the wide method and an independent sample of 22 other runners to use the narrow method, obtaining the same 44 times as in the actual study).



# First Base

Using an applet which tests a difference between these two means, ignoring the fact that it is paired data, we get a p-value of 0.3470.

This p-value is much larger than the one we obtained earlier.



Count samples:

Count = 347/1000 (0.3470)

# Theory-based Approach for Analyzing Data from Paired Samples, and M&Ms.

Section 7.3



# How Many M&Ms Would You Like?

Example 7.3

# How Many M&Ms Would You Like?

- Does your bowl size affect how much you eat?
- Brian Wansink studied this question with college students over several days.
- At one session, the 17 participants were assigned to receive either a small bowl or a large bowl and were allowed to take as many M&Ms as they would like.
- At the following session, the bowl sizes were switched for each participant.

# How Many M&Ms Would You Like?

- What are the observational units?
- What is the explanatory variable?
- What is the response variable?
- Is this an experiment or an observational study?
- Will the resulting data be paired?

# How Many M&Ms Would You Like?

The hypotheses:

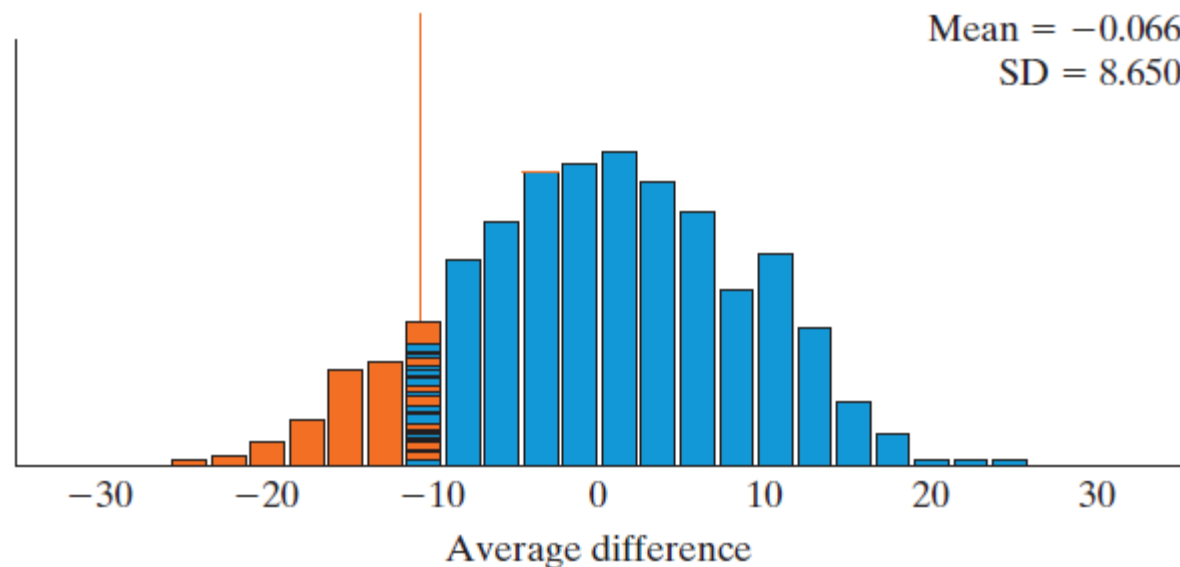
- $H_0: \mu_d = 0$ 
  - The long-run mean difference in number of M&Ms taken (small – large) is 0.
- $H_a: \mu_d < 0$ 
  - The long-run mean difference in number of M&Ms taken (small – large) is less than 0.

**TABLE 7.5** Summary statistics, including the difference (small – large) in the number of M&Ms taken between the two bowl sizes

Bowl size	Sample size, $n$	Sample mean	Sample SD
Small	17	$\bar{x}_s = 38.59$	$s_s = 16.90$
Large	17	$\bar{x}_l = 49.47$	$s_l = 27.21$
Difference = small – large	17	$\bar{x}_d = -10.88$	$s_d = 36.30$

# How Many M&Ms Would You Like?

- Here are the results of a simulation-based test.
- The p-value is quite large at 0.1220.

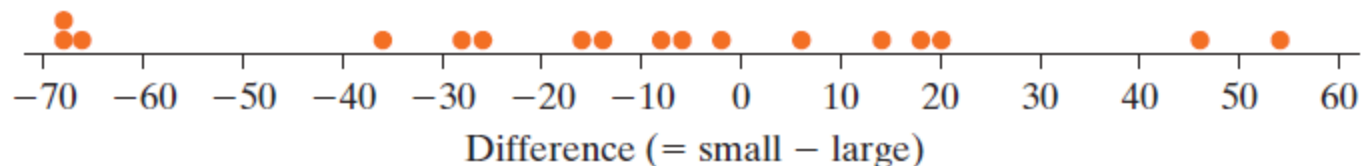


Count samples:

Count = 122/1000 (0.1220)

# How Many M&Ms Would You Like?

- Our sample size was only 17, but this distribution of differences is fairly symmetric and looks perhaps reasonably approximated by the normal dist., so we will proceed with a theory-based test here. It is a close call though.



# t-test.

- If we can assume the differences for each person are iid and normal with unknown sd, then with a theory based test we calculate the t-statistic:

$$t = \frac{\bar{x}_d}{s_d / \sqrt{n}}$$

- This kind of test is called a paired *t*-test.

# Conclusion

- The theory-based t-test results in a p-value of 11.72% and a 95% CI of (-29.5, 7.8). Thus it gives slightly different results than simulation, but we come to the same conclusion. We do not have strong evidence that the bowl size affects the number of M&Ms taken.
- We can see this in the large p-value (0.1172) and the confidence interval that includes zero (-29.5, 7.8).
- The confidence interval tells us that we are 95% confident that when given a small bowl, people will take somewhere between 29.5 fewer M&Ms to 7.8 more M&Ms on average than when given a large bowl.



# Why wasn't the difference statistically significant?

- There could be a number of reasons we didn't get significant results.
  - Maybe bowl size does not matter.
  - Maybe bowl size does matter and the difference was too small to detect with our small sample size.
  - Maybe bowl size does matter with some foods, like pasta or cereal, but not with a snack food like M&Ms.

# When to use which formula.

a. 1 sample numerical data, iid observations, want a 95% CI for  $\mu$ .

- If  $n$  is large and  $\sigma$  is known, use  $\bar{x} \pm 1.96 \sigma/\sqrt{n}$ .
- If  $n$  is small, draws are normal, and  $\sigma$  is known, use  $\bar{x} \pm 1.96 \sigma/\sqrt{n}$ .
- If  $n$  is small, draws are normal, and  $\sigma$  is unknown, use  $\bar{x} \pm t_{\text{mult}} s/\sqrt{n}$ .
- If  $n$  is large and  $\sigma$  is unknown,  $t_{\text{mult}} \sim 1.96$ , so we can use  $\bar{x} \pm 1.96 s/\sqrt{n}$ .

$n \geq 30$  is often considered large enough to use 1.96.

In practice, we typically do not know the draws are normal, but if the distribution looks roughly symmetrical without enormous outliers, the  $t$  formula may be reasonable.

b. 1 sample binary data, iid observations, want a 95% CI for  $\pi$ .

View the data as 0 or 1, so sample percentage  $\hat{p} = \bar{x}$ , and  $s = \sqrt{\hat{p}(1-\hat{p})}$ ,  $\sigma = \sqrt{\pi(1-\pi)}$ .

# When to use which formula.

a. 1 sample numerical data, iid observations, want a 95% CI for  $\mu$ .

- If  $n$  is large and  $\sigma$  is known, use  $\bar{x} \pm 1.96 \sigma/\sqrt{n}$ .
- If  $n$  is small, draws are normal, and  $\sigma$  is known, use  $\bar{x} \pm 1.96 \sigma/\sqrt{n}$ .
- If  $n$  is small, draws  $\sim$  normal, and  $\sigma$  is unknown, use  $\bar{x} \pm t_{\text{mult}} s/\sqrt{n}$ .
- If  $n$  is large and  $\sigma$  is unknown,  $t_{\text{mult}} \sim 1.96$ , so we can use  $\bar{x} \pm 1.96 s/\sqrt{n}$ .

b. 1 sample binary data, iid observations, want a 95% CI for  $\pi$ .

View the data as 0 or 1, so sample percentage  $\hat{p} = \bar{x}$ , and  
 $s = \sqrt{\hat{p}(1-\hat{p})}$ ,  $\sigma = \sqrt{\pi(1-\pi)}$ .

If  $n$  is large and  $\pi$  is unknown, use  $\bar{x} \pm 1.96 s/\sqrt{n}$ .

Here large  $n$  means  $\geq 10$  of each type in the sample.

What if  $n$  is small and the draws are not normal?

Then simulations are basically your only choice. There are other possible solutions outside the scope of this course, such as the bootstrap, which are sometimes useful in these situations.

# When to use which formula.

c. Numerical data from 2 samples, iid observations, want a 95% CI for  $\mu_1 - \mu_2$ .

If  $n$  is large and  $\sigma$  is unknown, use  $\bar{x}_1 - \bar{x}_2 \pm 1.96 \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$ .

As with one sample, if  $\sigma_1$  is known, replace  $s_1$  with  $\sigma_1$ , and the same for  $\sigma_2$ . And as with one sample, if  $\sigma_1$  and  $\sigma_2$  are unknown, the sample sizes are small, and the distributions are roughly normal, then use  $t_{\text{mult}}$  instead of 1.96. If the sample sizes are small, the distributions are normal, and  $\sigma_1$  and  $\sigma_2$  are known, then use 1.96.

d. Binary data from 2 samples, iid observations, want a 95% CI for  $\pi_1 - \pi_2$ .

Same as in c above, with  $\hat{p}_1 = \bar{x}_1$ ,  $s_1 = \sqrt{\hat{p}_1(1-\hat{p}_1)}$ ,  $\sigma_1 = \sqrt{\pi_1(1-\pi_1)}$ .

Large for binary data means sample has  $\geq 10$  of each type.

# When to use which formula.

e. Matched pairs data, iid observations, want a 95% CI for  $\mu$ .

Look at differences (score with treatment minus score with control) and treat differences as ordinary numerical data according to parts a or b.

- If  $n$  is large and  $\sigma$  is known, use  $\bar{x} \pm 1.96 \sigma/\sqrt{n}$ .
- If  $n$  is small, draws are normal, and  $\sigma$  is known, use  $\bar{x} \pm 1.96 \sigma/\sqrt{n}$ .
- If  $n$  is small, draws are normal, and  $\sigma$  is unknown, use  $\bar{x} \pm t_{\text{mult}} s/\sqrt{n}$ .
- If  $n$  is large and  $\sigma$  is unknown,  $t_{\text{mult}} \sim 1.96$ , so we can use  $\bar{x} \pm 1.96 s/\sqrt{n}$ .

$n \geq 30$  is often considered large enough to use 1.96.

In practice, we typically do not know the draws are normal, but if the distribution looks roughly symmetrical without enormous outliers, the  $t$  formula may be reasonable. This is often a tough judgement call.