# Lecture 20

## Hypothesis Testing

# Hypothesis Testing

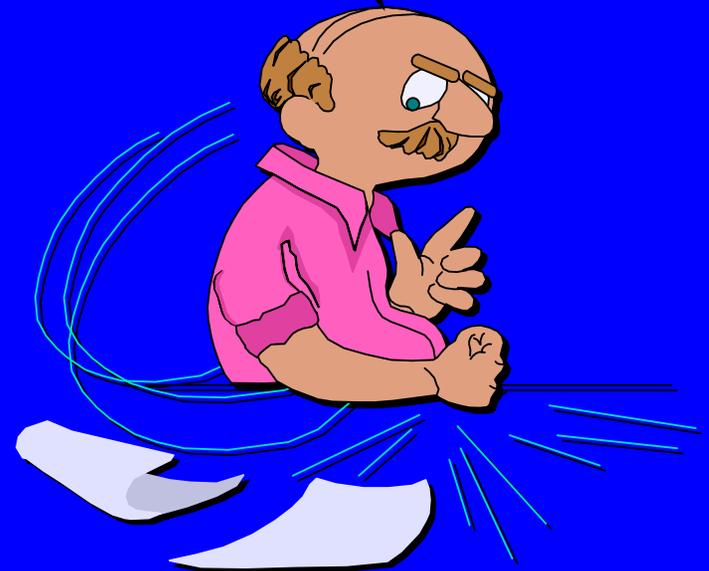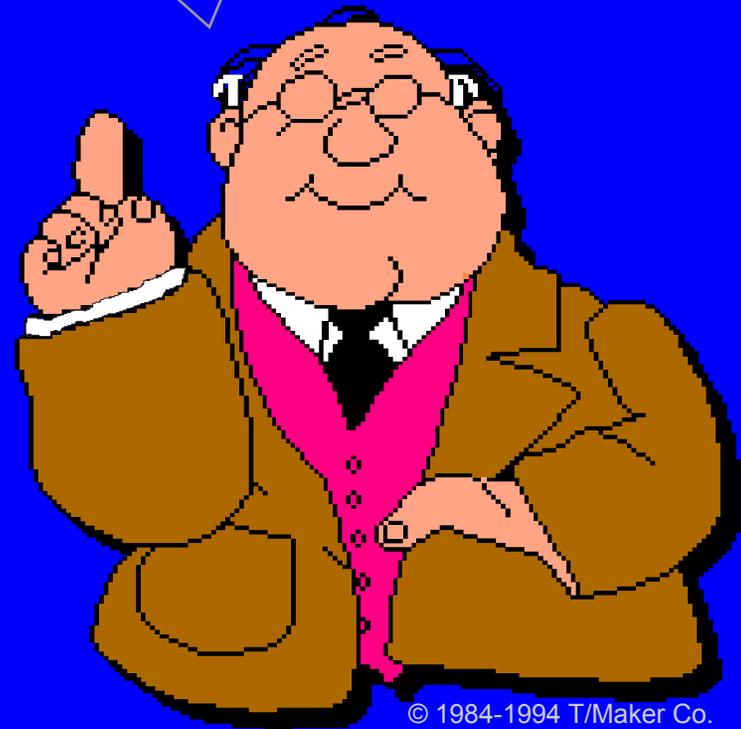# What's a Hypothesis?

A Belief about a Population Parameter

- Parameter Is **Population** Mean, Proportion, Variance

- Must Be Stated **Before** Analysis

I believe the mean GPA of this class is 3.5!

© 1984-1994 T/Maker Co.

# Thought Question 1:

Recall study where difference in sample means for weight loss based on dieting only versus exercising only was 3.2 kg. Same study showed difference in *average amount of fat weight lost was 1.8 kg* and corresponding *standard error was 0.83 kg*. Suppose means are actually equal, so mean difference in fat lost for populations is actually zero. What is the **standardized score** corresponding to observed difference of 1.8 kg? Would you expect to see a standardized score that large or larger very often?

# Thought Question 2:

In journal article in Case Study 6.4, comparing IQs for children of smokers and nonsmokers, one of the statements made was, "*After control for confounding background variables, the average difference [in mean IQs] observed at 12 and 24 months was 2.59 points (95% CI: –3.03, 8.20; P = 0.37).*"       (Olds et al., 1994, p. 223)

Reported value of 0.37 is the *p*-value. What are the **null and alternative hypotheses** being tested?

# Thought Question 3:

In chi-squared tests for two categorical variables, introduced in Chapter 13, we were interested in whether a relationship observed in a sample reflected a real relationship in the population.

What are the **null and alternative hypotheses**?

# Thought Question 4:

In Chapter 13, we found a statistically significant relationship between smoking (yes or no) and time to pregnancy (one cycle or more than one cycle).

Explain what the **type 1 and type 2 errors** would be for this situation, and the consequences of making each type of error.

# 23.1 How Hypothesis Tests Are Reported in the News

1. Determine the **null** hypothesis and the **alternative** hypothesis.
2. Collect and summarize the **data** into a **test statistic**.
3. Use the test statistic to determine the *p*-value.
4. The result is **statistically significant** if the *p*-value is less than or equal to the level of significance.

Often media only presents results of step 4.

# Example 1: Cranberry Juice

*CHICAGO (AP)  A scientific study has proven what many women have long suspected:* **Cranberry juice helps protect against bladder infections**. *Researchers found that* **elderly women who drank 10 ounces of a juice drink containing cranberry juice each day had less than half as many urinary tract infections as those who consumed a look-alike drink without cranberry juice**. *The study, which appeared today in the Journal of the American Medical Association, was funded by Ocean Spray Cranberries, Inc., but the company had no role in the study's design, analysis or interpretation, JAMA said. "This is the first demonstration that cranberry juice can reduce the presence of bacteria in the urine in humans," said lead researcher Dr. Jerry Avorn, a specialist in medication for the elderly at Harvard Medical School.*

<div align="right">

(*Davis Enterprise*, Mar. 9, 1994, p. A9)

</div>

# Example 1: Cranberry Juice

- Study to compare odds of getting an infection for population of elderly women following two regimes: 10 ounces of cranberry juice per day or 10 ounces of a placebo drink.

- **Null hypothesis** is odds ratio (juice/placebo) is 1.

- **Alternative hypothesis** is odds of infection are higher for the group drinking the placebo (juice/placebo < 1).

- The article indicates that the **odds ratio is under 50%**.

- Original article (Avorn et al., 1994) sets it at 42% and reports that the associated **$p$-value is 0.004**.

- Newspaper article captured most important aspect of research, but not clear that $p$-value was extremely low.

# 23.2 Testing Hypotheses About Proportions and Means

If the null and alternative hypotheses are expressed in terms of a **population proportion, mean, or difference between two means** and if the sample sizes are large …

… the **test statistic** is simply the corresponding **standardized score** computed assuming the null hypothesis is true; and the *p*-value is found from a table of percentiles for standardized scores.

**Step 1: Determine the null and alternative hypotheses.**

*Null hypothesis*: There is no difference in average fat lost in the population for the two methods.  The population mean difference is zero.

*Alternative hypothesis:* There is a difference in average fat lost in the population for the two methods.  The population difference is not zero.

$$H_0: m_d\text{-}m_e = 0$$

$$H_a: m_d\text{-}m_e \neq 0$$

When the alternative hypothesis includes a possible difference in either direction, the test is called a *two-sided*, or *two-tailed, hypothesis test*.  Our p-value needs to account for this.

**Step 2: Collect and summarize the data into a test statistic.**

The test statistic is the standardized score for the sample value when the null hypothesis is true.

If there is no difference in the two methods, then the mean *population* difference is zero.

The *sample* value is the observed difference in the two sample means:

$$\hat{m}_d - \hat{m}_e = 5.9 - 4.1 = 1.8 \, \text{kg}.$$

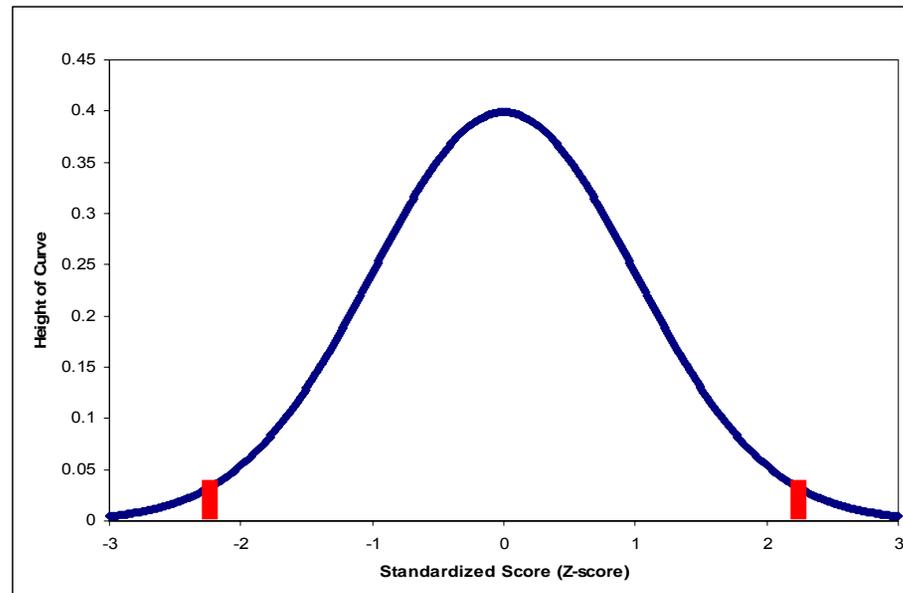The measure of uncertainty, or standard error for this difference is 0.83.

standardized score = z-score = (1.8 – 0)/0.83=2.17

**Step 3: Use the test statistic to determine the p-value.**

We need to consider a difference in both directions.

A standardized score of -2.17 would have been equally informative against the null hypothesis.

Using Table 8.1, the probability of a z-score of 2.17 or greater is 1 – 0.985=0.015.  The probability of a z-score of -2.17 or less is 0.015.  So our p-value is 2(0.015)=0.030.

**Step 4: Decide whether the result is statistically significant based on the p-value.**

If there were really no difference between dieting and exercise methods, we would see such an extreme result only 3% of the time.

We prefer to believe that the truth does not lie with the null hypothesis.

We conclude that there is a statistically significant difference between average fat loss for the two methods.

We reject the null hypothesis and accept the alternative hypothesis.

# Example 2: Public Opinion About President

On May 16, 1994, Newsweek reported the results of a public opinion poll that asked: *"From everything you know about Bill Clinton, does he have the honesty and integrity you expect in a president?"* (p. 23).

Poll **surveyed 518 adults and 233, or 0.45** of them (clearly less than half), answered yes.

Could Clinton's adversaries conclude from this that **only a minority (less than half) of the population** of Americans thought Clinton had the honesty and integrity to be president?

# Example 2: Public Opinion About President

## Step 1. Determine the null and alternative hypotheses.

*Null hypothesis:* There is no clear winning opinion on this issue; the proportions who would answer yes or no are each 0.50.

*Alternative hypothesis:* Fewer than 0.50, or 50%, of the population would answer yes to this question. The majority do not think Clinton has the honesty and integrity to be president.

## Step 2. Collect and summarize data into a test statistic.

Sample proportion is: $233/518 = 0.45$.

The *standard deviation* $= \sqrt{\dfrac{(0.50) \times (1 - 0.50)}{518}} = 0.022$.

*Test statistic:* $z = (0.45 - 0.50)/0.022 = -2.27$

# Example 2: Public Opinion About President

## Step 3. Determine the *p*-value.

*Recall the alternative hypothesis was one-sided.*
*p*-value = proportion of bell-shaped curve below –2.27
Exact *p*-value = 0.0116.

## Step 4. Make a decision.

The ***p*-value of 0.0116 is less than 0.05**, so we conclude that the proportion of American adults in 1994 who believed Bill Clinton had the honesty and integrity they expected in a president was **significantly less** than a majority.

**Step 1: Determine the null and alternative hypotheses.**

Null hypothesis: The proportion who would vote for Bush in the next election is half or greater.

Alternative hypothesis: Fewer than half of the population would vote for Bush in the next election.

$H_0$: p≥0.5

$H_a$: p<0.5

When the alternative hypothesis includes values in one direction only, the test is called a one-sided, or one-tailed, hypothesis test.  The p-value is computed using only the values in the direction specified in the alternative hypothesis.

**Step 2: Collect and summarize the data into a test statistic.**

Using the same procedure as in Example #1, the test statistic is:

$$z = (0.45 - 0.5) / \sqrt{(0.5)(0.5)/518} = -2.27.$$

**Step 3: Use the test statistic to determine the p-value.**

The p-value is the probability of observing a standardized score of -2.27 or less just by chance.

Using Table 8.1 (or Excel), this value is 0.0116.

**Step 4: Decide whether the result is statistically significant based on the p-value.**

Using the 0.05 criterion, we have found a statistically significant result.

Therefore we reject the null hypothesis and accept the alternative. We believe that fewer than half of all Americans plan to vote for Bush in the next election.

**Step 1: Determine the null and alternative hypotheses.**

Null hypothesis: In the population of young drivers, there is no relationship between gender and whether the driver drank alcohol in the last 2 hours.

Alternative hypothesis: In the population of young drivers, one of the two sexes is more likely than the other to have consumed alcohol in the 2 hours prior to driving.

**Step 2: Collect and summarize the data into a test statistic.**

The Chi-squared test statistic for this data is 1.637.

**Step 3: Use the test statistic to determine the p-value.**

The p-value corresponding to this data is 0.2007.

**Step 4: Decide whether the result is statistically significant based on the p-value.**

The p-value is 0.20 (> 0.05), so we would probably not be willing to rule out chance as an explanation for the observed difference in the proportions.

We fail to reject the null hypothesis.

# 23.3 Revisiting Case Studies: How Journals Present Tests

Whereas newspapers and magazines tend to simply report the decision from hypothesis testing, **journals tend to report $p$-values as well**.

This allows you to make your own decision, based on the severity of a type 1 error and the magnitude of the $p$-value.

# Case Study 6.1: Mozart, Relaxation, and Performance on Spatial Tasks

**Three listening conditions**— Mozart, a relaxation tape, and silence—and all subjects participated in all three conditions.

*Null hypothesis:* **No differences** in population mean spatial reasoning IQ scores after each of three listening conditions.

*Alternative hypothesis:* Population mean spatial reasoning IQ scores **do differ for at least one** of the conditions compared with the others.

# Case Study 6.1: Mozart, Relaxation, and Performance on Spatial Tasks

*A one-factor (listening condition) repeated measures analysis of variance ... revealed that subjects performed better on the abstract/spatial reasoning tests after listening to Mozart than after listening to either the relaxation tape or to nothing (F[2,35] = 7.08, **p = 0.002**).*

*(Rauscher et al., 14 October 1993, p. 611)*

**Conclusion:** At least one of the means differs from the others. If there were no population differences, sample mean results would vary as much as the ones in this sample did, or more, only 2 times in 1000 (0.002).

# Case Study 6.1: Mozart, Relaxation, and Performance on Spatial Tasks

*The music condition differed significantly from both the relaxation and silence conditions (Scheffé's t = 3.41, p = 0.002; t = 3.67, p = 0.0008, two-tailed, respectively). The relaxation and silence conditions did not differ (t = 0.795, p = 0.432, two-tailed).*

*(Rauscher et al., 14 October 1993, p. 611)*

Significant differences were found between the music and relaxation conditions (*p*-value = 0.002) and between the music and silence conditions (*p*-value = 0.0008). The difference between the relaxation and silence conditions, however, was not statistically significant (*p*-value = 0.432).

# Case Study 5.1: Quitting Smoking with Nicotine Patches

Compared the smoking cessation rates for smokers randomly assigned to use a nicotine patch versus a placebo patch.

*Null hypothesis:* The proportion of smokers in the population who would quit smoking using a nicotine patch and a placebo patch are the **same**.

*Alternative hypothesis:* The proportion of smokers in the population who would quit smoking using a **nicotine patch is higher** than the proportion who would quit using a placebo patch.

# Case Study 5.1:  Quitting Smoking with Nicotine Patches

*Higher smoking cessation rates were observed in the active nicotine patch group at 8 weeks (46.7% vs 20%) (P < .001) and at 1 year (27.5% vs 14.2%) (P = .011).*

*(Hurt et al., 1994, p. 595)*

**Conclusion:**  *p*-values are quite small: less than 0.001 for difference after 8 weeks and equal to 0.011 for difference after a year. Therefore, rates of quitting are significantly higher using a nicotine patch than using a placebo patch after 8 weeks and after 1 year.

# Case Study 6.4: Smoking During Pregnancy and Child's IQ

Study investigated impact of maternal smoking on subsequent IQ of child at ages 1, 2, 3, and 4 years of age.

*Null hypothesis:* Mean IQ scores for children whose mothers smoke 10 or more cigarettes a day during pregnancy **are same** as mean for those whose mothers do not smoke, in populations similar to one from which this sample was drawn.

*Alternative hypothesis:* Mean IQ scores for children whose mothers smoke 10 or more cigarettes a day during pregnancy are **not the same** as mean for those whose mothers do not smoke, in populations similar to one from which this sample was drawn.

# Case Study 6.4: Smoking During Pregnancy and Child's IQ

*Children born to women who smoked 10+ cigarettes per day during pregnancy had developmental quotients at 12 and 24 months of age that were 6.97 points lower (averaged across these two time points) than children born to women who did not smoke during pregnancy (95% CI: 1.62,12.31, P = .01); at 36 and 48 months they were 9.44 points lower (95% CI: 4.52, 14.35, P = .0002).* (Olds et al., 1994, p. 223)

Researchers conducted *two-tailed tests* for possibility the mean IQ score could actually be *higher* for those whose mothers smoke. The CI provides evidence of the direction in which the difference falls. The *p*-value simply tells us there is a statistically significant difference.

# For Those Who Like Formulas

**Some Notation for Hypothesis Tests**

The null hypothesis is denoted by $H_0$, and the alternative hypothesis is denoted by $H_1$ or $H_a$.

"alpha" $= \alpha =$ desired probability of making a type 1 error when $H_0$ is true; we reject $H_0$ if $p$-value $\leq \alpha$.

"beta" $= \beta =$ probability of making a type 2 error when $H_1$ is true; power $= 1 - \beta$

**Steps for Testing the Mean of a Single Population**

Denote the population mean by $\mu$ and the sample mean and standard deviation by $\overline{X}$ and $s$, respectively.

**Step 1.** $H_0: \mu = \mu_0$, where $\mu_0$ is the *chance* or *status quo* value.

$H_1: \mu \neq \mu_0$ for a two-sided test; $H_1: \mu < \mu_0$ or $H_1: \mu > \mu_0$ for a one-sided test, with the direction determined by the research hypothesis of interest.

**Step 2.** This test statistic applies only if the sample is large. The test statistic is

$$z = \frac{\overline{X} - \mu_0}{s / \sqrt{n}}$$

# For Those Who Like Formulas

**Step 3.** The $p$-value depends on the form of $H_1$. In each case, we refer to the proportion of the standard normal curve above (or below) a value as the "area" above (or below) that value. Then we list the $p$-values as follows:

| Alternative Hypothesis | $p$-Value |
| --- | --- |
| $H_1: \mu \neq \mu_0$ | $2 \times$ area above $|z|$ |
| $H_1: \mu > \mu_0$ | area above $z$ |
| $H_1: \mu < \mu_0$ | area below $z$ |

**Step 4.** You must specify the desired $\alpha$; it is commonly 0.05. Reject $H_0$ if $p$-value $\leq \alpha$.