

# Hypothesis Tests for $\mu$

Making Decisions About the True  
Population Mean

# Confidence Intervals vs Hypothesis Tests

- Both confidence and hypothesis tests are used to make inferences and decisions about the population parameters of interest based on samples. They accomplish the same goals.
- Tests or confidence intervals are used more or less in certain subject areas. A general tendency is that CI's are used more in the sciences and tests more in the Social Sciences. However, CI's are becoming more common in the soc sciences.

# Hypothesis Test Demonstration

- The terminology for hypothesis tests is very thick. This demonstration should help you remember what we are doing when things get tough.
- I will flip a coin several times. Keep track of the results.

# Demonstration Result

- As I got too many tails in a row, about 4 or 5 usually, most of you started to doubt the fairness of the flipping system.
- When we got way too many tails in a row, you knew there was a problem.
- The data got too unusual or atypical to be coming from a fair 50-50 system. So you rejected your notion of fairness and claimed you were being tricked.

# Demonstration Result

- Just remember that when data gets too unusual you will doubt your underlying hypothesis about how the world works, just like in our demonstration.
- $P(5 \text{ tails in a row}) = (1/2)^5 = .031$ , this is about when most people start to object.
- Keeping this example in your mind will help understand things as we continue.

# Four Stages of a Hypothesis Test

- 1) Write hypotheses: null hypothesis  $H_0$ , and alternative hypothesis,  $H_a$ .
- 2) Calculate the test statistic.
- 3) Calculate the p-value.
- 4) Give a conclusion.

# Null Hypothesis

- The null hypothesis is denoted  $H_0$ .
- This is pronounced “H – not”. If you call it “Hoe”, I will chuckle at you because it sounds dorky.
- The null hypothesis is a belief about the world, or population. It is typically what we are trying to show is false.

# Alternative Hypothesis

- The alternative hypothesis is denoted  $H_a$ .
- This hypothesis is pronounced “H-A” or alternative hypothesis. If you say “Hah”, I will chuckle because you will sound dorky.
- The alternative hypothesis is a belief about the population that you believe is true and are probably trying to show is true.



# Test Statistic

- The test statistic is a convenient summary of the sample data that can be easily used to make decisions about the hypotheses.
- It is often some kind of transform of the data to a more convenient form, like Z-scores.
- Test statistic is calculated under assumption the  $H_0$  is true.

# P-Value

- The next step in a test is the p-value.
- A p-value is a probability, so it is a number that is between 0 and 1.
- It is a measure of how consistent the sample data is with the null hypothesis.
- There are two definitions of p-value that will be useful to you:

# Definitions

- The p-value is the probability of observing a value of the test statistic as extreme or more extreme than the one observed, if the  $H_0$  is true.
- The p-value is the probability of observing data like ours if the null hypothesis is true.

# Conclusion

- The conclusion gives our final opinion about what the data tells us about our hypotheses about the world. A good conclusion will contain:
  - 1) Data consistency/inconsistency with  $H_0$ ?
  - 2) Evidence to doubt  $H_0$ ?
  - 3) Evidence to support  $H_a$ ?
  - 4) Final conclusion about problem, final punch line. (Visit Eiffel Tower!)

# An Example

- Tobacco company claims new cigarette is low-tar, less than 5mg of tar per cigarette on average. We work for consumer organization that tests such claims. From our perspective the hypotheses are:
- $H_0: \mu \leq 5.0\text{mg}$      $H_a: \mu > 5.0\text{mg}$
- This is because if we believe  $\mu$  is above 5mg, we will need to take some action against the tobacco company.

# Tar Example

- The data for this experiment was a sample mean of  $\bar{x} = 5.5$  mg, and a standard deviation of  $\sigma = 1.2$  mg/cig. There were  $n = 36$  cigs tested.
- Test statistic is calculated assuming the  $H_0$  is true. We use a  $\mu$  value from the null  $H_0$ . The test statistic value is  $Z =$

$$Z = \frac{\bar{x} - \mu_0}{\frac{\sigma}{\sqrt{n}}} = \frac{5.5 - 5.0}{1.2 / \sqrt{36}} = \frac{.5}{.2} = 2.5$$

# Tar P-Value

- The p-value is the probability of observing a value of  $Z$  more extreme than 2.5 if the  $H_0$  was true.
- What is more extreme comes from the  $H_a$  hypothesis, so that our p-value is:
- $P(Z > 2.5)$  because our  $H_a$  had a greater than ( $>$ ) sign.
- From table, p-value = .0062.

# Conclusion

- Let's do the recipe for the conclusion:
- 1) The data are unusual and unlikely to occur if the  $H_0$  was true, so data are inconsistent with  $H_0$ .
- 2) There is evidence to doubt  $H_0$ .
- 3) There is evidence to support  $H_a$ . (Must be !)
- 4) There is evidence that the mean  $\mu$  is above 5mg, so evidence against the company claim, and this is evidently not a low-tar cigarette.
- DIAL-LAWYERS !!



# Example

- Page 439 in text says Tim Kelley's weight on driver license is 187, but we suspect he has gained weight since last license. This means we have:
- $H_0: \mu = 187$  versus  $H_a: \mu > 187$
- Sample data weights from this month: 190.5, 189, 195.5, 187 give  $\bar{x}$  190.5,  $n=4$ , and assume  $\sigma=3$ , and wts are normally distributed.

# Example

- Test Statistic:

$$Z = \frac{\bar{x} - \mu_0}{\frac{\sigma}{\sqrt{n}}} = \frac{190.5 - 187}{3 / \sqrt{4}} = \frac{3.5}{1.5} = 2.33$$

# Example: P-Value

- P-value is the probability of observing a test statistic  $Z$  value more extreme than 2.33. Because the  $H_a$  was greater than, we need to find the chance of a test statistic greater than 2.33,  $P\text{-value} = P(Z > 2.33) = .0099$ .

# Conclusion

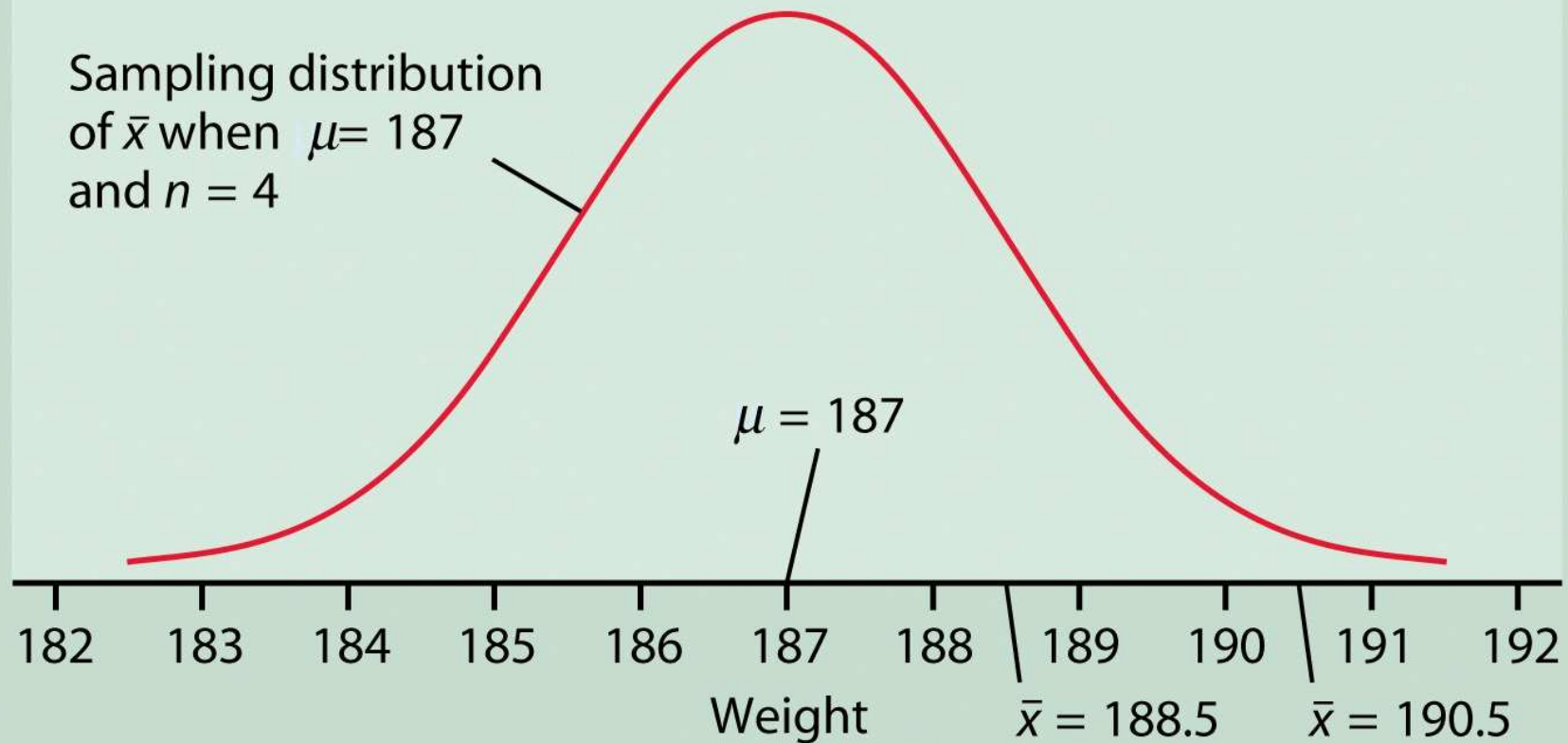
- 1) The data are unlikely to occur if the  $H_0$  was true, the data are inconsistent with  $H_0$ .
- 2) There is evidence to doubt the  $H_0$ .
- 3) There is evidence to support the  $H_a$ .
- 4) There is evidence that the average weight  $\mu$ , is above 187 pounds, and evidence that Tim Kelley has gained weight since the last license.

# Conclusion Note

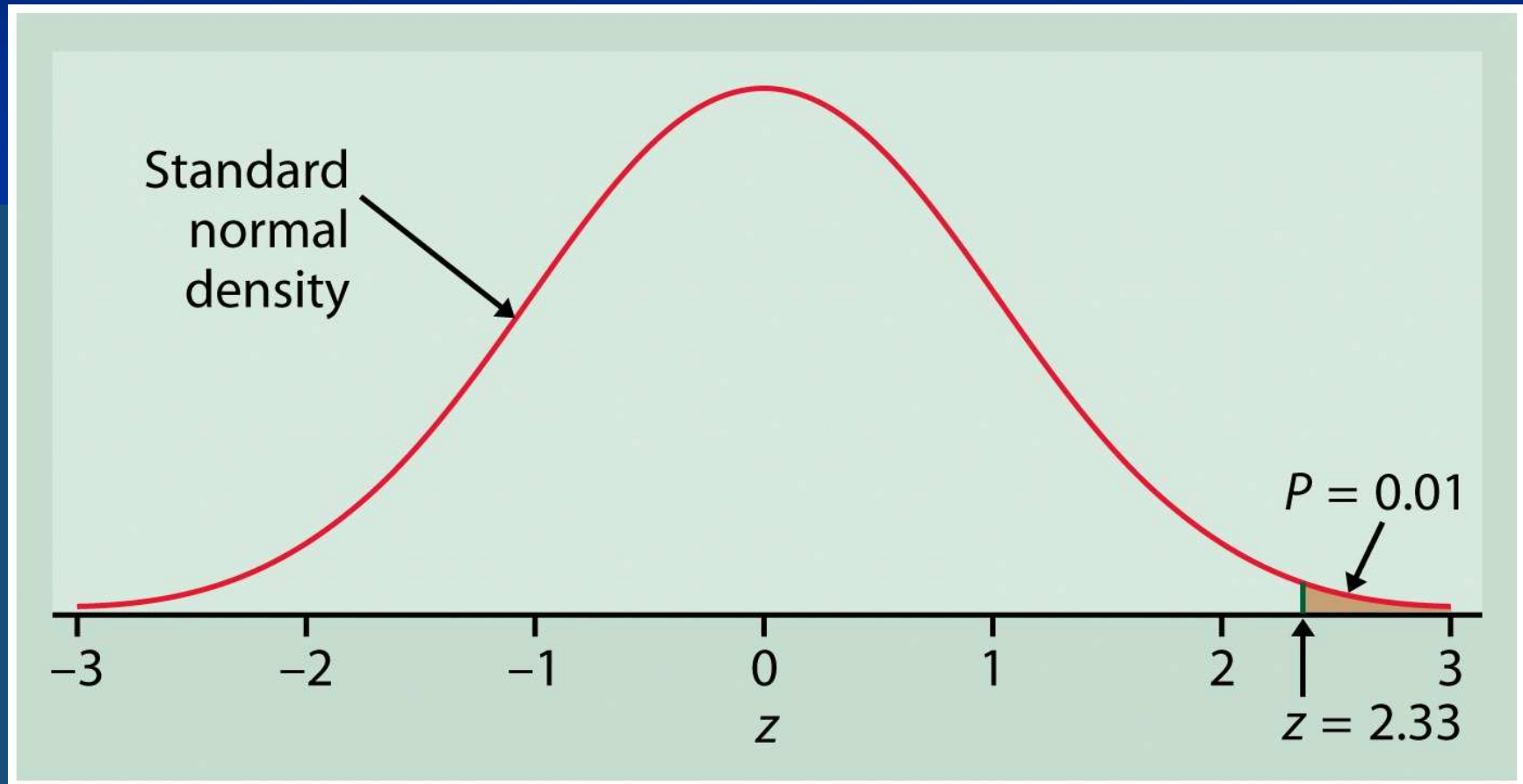
- Notice that the first steps are full of jargon, but that the end statement should be clearly understood by anyone reading it. It simply gives the plain English answer to the research question.
- Notice that these examples are just like the coin flipping demonstration in this unit – the data were weird and contradictory to the  $H_0$ , and so we begin to doubt the  $H_0$  and support  $H_a$ .

# Weight Example

Sampling distribution  
of  $\bar{x}$  when  $\mu = 187$   
and  $n = 4$



# Weight Example



# Psychology Today Example

- Bohrnstedt-Knoke example. Researchers surveyed 2013 readers of Psychology Today and asked them to rate their physical attractiveness to others. Scale: 1 = much less attractive than others, 4= About the same as others, to 7= much more attractive than others.
- Rate your own physical attractiveness on scale of 1 to 7. Turn in slips of paper, NO NAMES PLEASE !



# Psychology Today

- Researchers hypothesized that people will rate their own physical attractiveness above average.
- Research hypothesis ( $H_a$ ): people tend to over-rate their own physical attractiveness.
- This means null hypothesis is that people do not rate themselves above average.
- $H_o: \mu \leq 4$ ,  $H_a: \mu > 4$ .

# Psychology Today Example

- Data:  $\bar{x} = 4.9$ ,  $\sigma = 1.153$ ,  $n = 2013$ .
- Test Statistic:

$$Z = \frac{\bar{x} - \mu_0}{\sigma / \sqrt{n}} = \frac{4.9 - 4.0}{1.153 / \sqrt{2013}} = 35.02$$

# Psychology Today

- P-value =  $P(Z > 35.02) = \text{about } 0$ .
- It is basically impossible to get out past 35.02 on the normal scale. Way off the chart !

# Psychology Today

- Conclusion:
- 1) Data is unlikely (impossible) to occur if the  $H_0$  was true, data inconsistent with the  $H_0$ .
- 2) Lots of evidence to doubt the  $H_0$ .
- 3) Evidence to support the  $H_a$ .
- 4) Evidence that mean rating is above 4 which means evidence that average rating is above average. So evidence people do over-rate themselves, on average. Why do you think?

# Stat 1601 Student Self - Ratings

- Class data:  $\bar{x} = 4.41$ ,  $\sigma = 1.153$ ,  $n=79$
- Test Statistic:  $Z = 3.133$
- P-Value = .0009
- Evidence stat students rate themselves above average.
- Note that stat students are all above average, and physically very hot, Hot, HOT !
- This is well known on campus. Ask anyone.

# Stat 1601 Results

## Sample Statistics for rating

N	Mean	Std. Dev.	Std. Error
79	4.41	1.16	0.13

Null hypothesis: Mean of rating  $\leq 4$

Alternative: Mean of rating  $> 4$

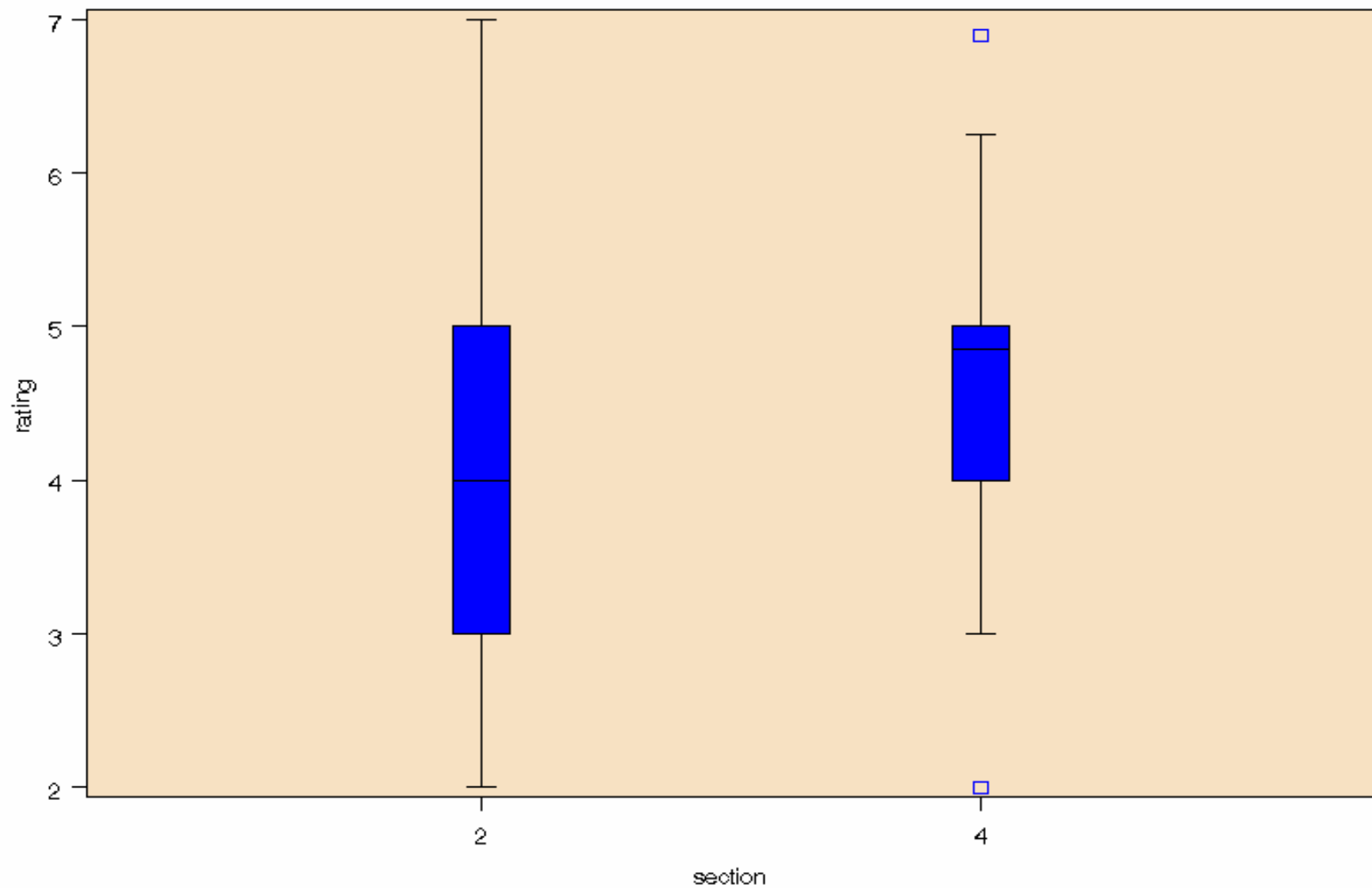
Z Statistic	Prob > Z
3.133	0.0009

## 95% Confidence Interval for the Mean

Lower Limit	Upper Limit
4.15	4.66

# Stat 1601 Results

Physical Attractiveness



# Physical Attractiveness

- Notice that if we would use any other value for  $\mu_0$  in the null hypothesis other than 4, say 2.74, the test statistic becomes even bigger and we still end up doubting the  $H_0$ . No matter what we put for  $\mu_0$  we get the same conclusion.
- Shows that we only need to use values on the border between  $H_0$  and  $H_a$ .



# Hypothesis Test Example

- $H_0$ : 23 star papers, 1 plain
- $H_a$ : 1 star, 23 plain
- Data:
- $P\text{-Value} = P(1 \text{ plain} \mid H_0) = 1/24 = .042$
- $P\text{-Value} = P(1 \text{ star} \mid H_0) = 23/24 = .96$
- $P\text{-Value} = P(2 \text{ plain} \mid H_0) = 0 / \# \text{ ways} = 0$
- A plain or two plains give lots of evidence to doubt the  $H_0$ .

# Shakespeare Example

- Video clip on new Bill Shakespeare poem.
- $H_0$ : Bill S. wrote poem,  $H_a$ : Bill S not write it
- Test Statistic is based on number of new words not used before in a Bill S poem. Value:  $Z = \text{approx } 1.0$ .
- P-Value = about .15
- Conclusion: 1) Data is consistent with  $H_0$ .
- 2) No evidence to doubt  $H_0$ .
- 3) No evidence to support  $H_a$ .
- 4) No evidence that Bill did not write poem, so no reason to doubt that he wrote it. He might not have written it, but we have no evidence to say he didn't ! Included in complete works.

# Stylometry Authorship

- Recently computer analysis techniques have been used to establish authorship of several disputed documents. An example is the Federalist papers. Although they were published anonymously, the author of 73 of these was determined to be John Jay (5) and the rest divided between Alexander Hamilton and James Madison. There were twelve that were left open to question. Using frequency of usage of the small filler words, they found overwhelming evidence favoring Madison as the author of all twelve disputed papers.
- A second example deals with an unfinished novel by Jane Austen when she died in 1817. A skilled author completed the novel and had it published. Although she duplicated the style she failed to duplicate the subconscious habits of detail. When these habit patterns were examined, the difference was clearly evident.

# Authorship Situations

- The Federalist Papers - who wrote the essays whose authorship is unknown?
- The Epistle to the Hebrews - was this written by St Paul (as traditionally believed?), or are modern scholars right to dispute this?
- The Acts of the Apostles - is there evidence to support the traditional ascription of this to the author of Luke's Gospel?
- Plato's Dialogues - can stylometric analysis shed light on their probable order of composition?
- Aristotle's Ethics - did the chapters which are common to the Nicomachean and the Eudemian Ethics originate as part of the Nicomachean Ethics, or as part of the Eudemian Ethics?
- The Shakespeare controversy - how plausible is the claim that the plays traditionally ascribed to Shakespeare were really written by Bacon (or perhaps by Marlowe)?

# Radon Detectors

- Manufacturer wants to determine if home radon detectors are reliable enough to sell. Detectors exposed to known source of 105 picocuries.
- Null hypothesis:  $H_0: \mu = 105$  versus  $H_a: \mu \neq 105$  (not equal to 105). Why not equal?
- If the actual average reading was over 105 it would mean the detectors were over-stating the radon risk, and we would freak out homeowners. Bad idea.

# Radon Detectors

- If the average reading was under 105, the detectors would be understating the risk and homeowners would think things were fine when in fact there was much radon risk. Ooops! Not Good!
- This is why the  $H_a$  is not equal to 105, we care about values of  $\mu$  too much or too little compared to 105.



# Radon Detectors

- Data:  $\bar{x}$  104.133,  $\sigma=9$ ,  $n=12$ .
- Test Statistic:  $Z = -.33$ .
- P-Value=Probability of observing  $Z$  value more extreme than  $-.33$ . More extreme is more negative than  $-.33$ , AND also more positive than  $.33$  because this would also be further from center of 0 than  $-.33$ . So...
- $P\text{-Value} = P(Z < -.33) + P(Z > .33) =$
- $= 2 * P(Z > |.33|) = 2(.3707) = \text{big chance}$

# Radon Conclusion

- Conclusion: 1) Data is likely to occur if  $H_0$  was true, so data is consistent with  $H_0$ .
- 2) No evidence to doubt  $H_0$ .
- 3) No evidence to support  $H_a$ .
- 4) No evidence to suggest the mean reading is anything other than 105picocuries. So no evidence of a problem with the detectors, so go ahead and sell them.



# Caution !

- A big p-value as in the last example (over .1 say), only means we have no reason to bad-mouth the  $H_0$ .
- It does NOT mean we have proved the  $H_0$  is true !!!
- The next example will hopefully illustrate the idea.

# OJ

- Consider the OJ trial, the  $H_0$  was  $H_0$ : OJ Innocent,  $H_a$ : OJ Guilty.
- Test Statistic = compilation of evidence
- P-value = big chance, over .1. Jury thought incriminating evidence such as this could easily occur with LA police investigating.
- Big p-value means jury did not have enough evidence to dump  $H_0$ , and support  $H_a$ , and convict OJ.

# OJ Continued

- $H_0$ : OJ Innocent ,  $H_a$ : OJ Guilty
- So a big p-value means not enough evidence to convict.
- It does not mean we proved that the  $H_0$  was true!!!  
Did the jury prove that OJ was not the killer?
- No, they simply weren't convinced he did it, they did not prove he didn't do it.
- A big p-value does not necessarily mean  $H_0$  is true, it might indeed be true, but can't conclude this.

# What is a small p-value?

- Most of scientific world uses .05 as the cut-off.
- P-values less than .05 are considered small chances and those above .05 are large.
- This gets ridiculous when p-values = .0501 or .049999. P-values are in shades of gray, not black and white.

# What I Expect

- If a p-value is below .05, you need to conclude it is a small chance and evidence to doubt  $H_0$ .
- If a p-value is above .1, you need to conclude it is a big chance and no evidence to doubt  $H_0$ .
- For p-values between .05 and .1 you have the option of concluding either the chance is big or small. I will accept a correctly argued position.

# Is mean body temp 98.6 degrees?

- $H_0: \mu = 98.6$  vs  $H_a: \mu \neq 98.6$  degrees
- $Z = (98.12 - 98.6) / (.63 / \sqrt{93}) = -7.35$
- $P\text{-Value} = P(Z < -7.35) + P(Z > 7.35) =$
- $= 2 * P(Z < -7.35) = 2 * \text{zero} = 0$
- Conclusion:

Is Mean Body Temp 98.6 degs?

$$H_0 : \mu = 98.6$$

$$H_a : \mu \neq 98.6$$

$$Z = \frac{98.12 - 98.6}{\frac{.63}{\sqrt{93}}} = -7.35$$

$$P\text{-Value} = 2 * P(Z < -7.35) = 0$$

# Body Temp Conclusion

- The data are very unusual if the  $H_0$  was true, the data are very inconsistent with the  $H_0$ .
- There is evidence to doubt the  $H_0$ .
- There is evidence to support the  $H_a$ .
- There is evidence the mean body temp is not 98.6 degrees. Because the sample mean,  $\bar{x}$  is less than 98.6, we have evidence the true mean body temp is below 98.6.



# Consequences of Incomplete Conclusion

■ Dear Unemployed,

Hit the bricks dude! What a pathetic answer to our question. We pay you statisticians big money. The least you could do is be kind enough to answer our questions. If we have evidence to doubt the  $H_0$ , fine, but which way does the evidence go you lunkhead!

This is the last time we hire anyone from Podunk State University. Next time we'll recruit hard and hire someone from an elite school like UMM where the stat professors teach the students how to present conclusions!

Sincerely,

Your Boss

PS. Would you like fries with that?

# Cancer Therapy Example

- Old treatment gives average remaining life of  $\mu=4.2$  years. New therapy experiment gives  $\bar{x}=4.5$  years,  $\sigma=1.1$  years,  $n=80$  patients.
- Believe new therapy prolongs life.
- $H_0: \mu=4.2$  vs  $H_a: \mu > 4.2$
- $Z = [4.5 - 4.2] / [1.1 / \sqrt{80}] = 2.44$
- $P\text{-Value} = P(Z > 2.44) = .0073$  (yippee)  $= .007$

# Cancer Conclusion

- 1) The data are unusual if  $H_0$  was true, data is inconsistent with the  $H_0$ .
- 2) There is evidence to doubt  $H_0$ .
- 3) There is evidence to support the  $H_a$ .
- 4) There is evidence the mean remaining life is above 4.2 years, and so evidence the new therapy has longer average life than old treatment.