

## UCLA STAT XL 10 Introduction to Statistical Reasoning

**Instructor: Ivo Dinov,**  
Asst. Prof. In Statistics and Neurology

University of California, Los Angeles, Spring 2002  
<http://www.stat.ucla.edu/~dinov/>

STAT XL 10, UCLA, Ivo Dinov

Slide 1

## Chapter 26: Significance Testing -- Using Data to Test Hypotheses

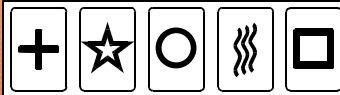
- Getting Started
- What do we test? Types of hypotheses
- Measuring the evidence against the null
- Hypothesis testing as decision making
- Why tests should be supplemented by intervals

STAT XL 10, UCLA, Ivo Dinov

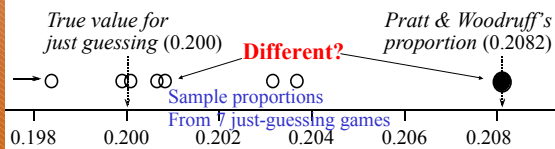
Slide 2

### ESP (extra sensory perception) or just guessing?

Deck of equal number of Zener/Rhine cards



n=60,000 random draws resulting in 12,489 correct guesses

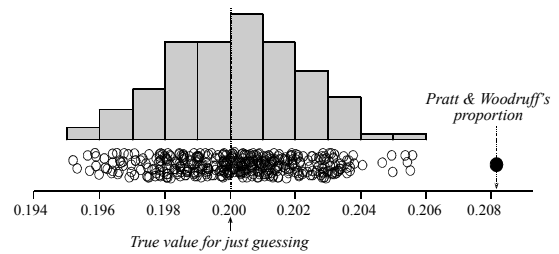


Can sampling variations alone account for Pratt & Woodruff's success rate = 20.82% correct vs. 20% expected.

Slide 3

STAT XL 10, UCLA, Ivo Dinov

### ESP or just guessing?



**Sample proportions from 400  
"just-guessing" experiments**

Slide 4

STAT XL 10, UCLA, Ivo Dinov

### Was Cavendish's experiment biased?

A number of famous early experiments of measuring physical constants have later been shown to be biased.

#### Mean density of the earth

True value = 5.517

**Cavendish's data:** (from previous Example)

5.36, 5.29, 5.58, 5.65, 5.57, 5.53, 5.62, 5.29, 5.44, 5.34, 5.79, 5.10, 5.27, 5.39, 5.42, 5.47, 5.63, 5.34, 5.46, 5.30, 5.75, 5.68, 5.85

n = 23, sample mean = 5.483, sample SD = 0.1904

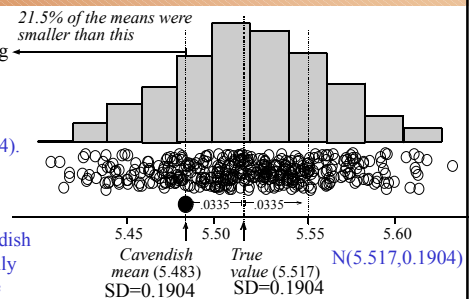
Slide 5

STAT XL 10, UCLA, Ivo Dinov

### Was Cavendish's experiment biased?

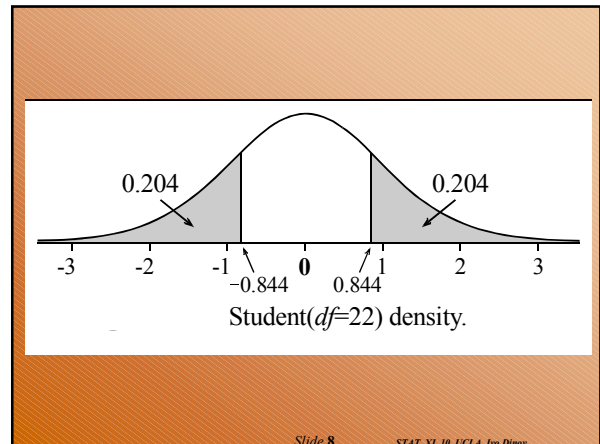
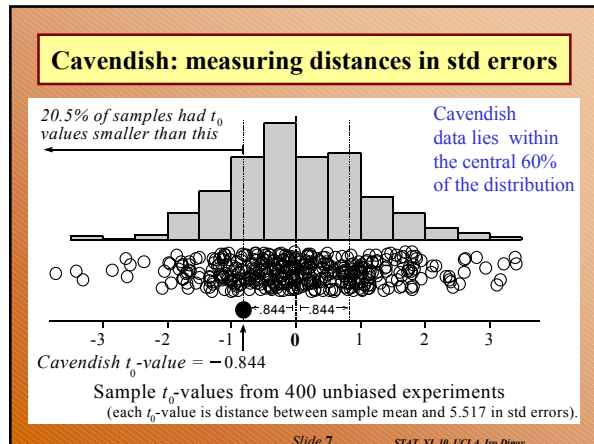
Simulate taking 400 sets of 23 measurements from N(5.517, 0.1904). Plotted are the results of the sample means.

Are the Cavendish values unusually diff. From true mean?



Slide 6

STAT XL 10, UCLA, Ivo Dinov



### Measuring the distance between the true-value and the estimate in terms of the SE

- Intuitive criterion: Estimate is credible if it's not far away from its hypothesized true-value!
- But how far is far-away?
- Compute the distance in standard-terms:

$$T = \frac{\text{Estimator} - \text{TrueParameterValue}}{\text{SE}}$$

- Reason is that the distribution of  $T$  is known in some cases (Student's  $t$ , or  $N(0,1)$ ). The estimator (obs-value) is typical/atypical if it is close to the center/tail of the distribution.

Slide 9 STAT.XL.10.UCLA.Im.Dinov

### Comparing CI's and significance tests

- These are different methods for coping with the uncertainty about the true value of a parameter caused by the sampling variation in estimates.
- **Confidence interval:** A fixed level of confidence is chosen. We determine a range of possible values for the parameter that are consistent with the data (at the chosen confidence level).
- **Significance test:** Only one possible value for the parameter, called the hypothesized value, is tested. We determine the strength of the evidence (confidence) provided by the data against the proposition that the hypothesized value is the true value.

Slide 10 STAT.XL.10.UCLA.Im.Dinov

### Review

- What intuitive criterion did we use to determine whether the hypothesized parameter value ( $p=0.2$  in the ESP Example, and  $\mu = 5.517$  in Earth density ex.) was credible in the light of the data? (Determine if the data-driven parameter estimate is consistent with the pattern of variation we'd expect get if hypothesis was true. If hypothesized value is correct, our estimate should not be far from its hypothesized true value.)
- Why was it that  $\mu = 5.517$  was credible in Ex. 2, whereas  $p=0.2$  was not credible in Ex. 1? (The first estimate is consistent, and the second one is not, with the pattern of variation of the hypothesized true process.)

Slide 11 STAT.XL.10.UCLA.Im.Dinov

### Review

- What do  $t_0$ -values tell us? (Our estimate is typical/atypical, consistent or inconsistent with our hypothesis.)
- What is the essential difference between the information provided by a confidence interval (CI) and by a significance test (ST)? (Both are uncertainty quantifiers. CI's use a fixed level of confidence to determine possible range of values. ST's one possible value is fixed and level of confidence is determined.)

Slide 12 STAT.XL.10.UCLA.Im.Dinov

## Hypotheses

**Guiding principles**

We cannot rule in a hypothesized value for a parameter, we *can only* determine whether there is evidence *to rule out* a hypothesized value.

The null hypothesis tested is typically a skeptical reaction to a *research hypothesis*

Slide 13 STAT XL 10, UCLA, Ivo Dinov

## Comments

- Why can't we (**rule-in**) prove that a hypothesized value of a parameter is exactly true? (Because when constructing estimates based on data, there's always sampling and may be non-sampling errors, which are normal, and will effect the resulting estimate. Even if we do 60,000 ESP tests, as we saw earlier, repeatedly we are likely to get estimates like 0.2 and 0.200001, and 0.199999, etc. – non of which may be exactly the theoretically correct, 0.2.)
- Why use the rule-out principle? (Since, we can't use the rule-in method, we try to find compelling evidence against the observed/data-constructed estimate – to reject it.)
- Why is the null hypothesis & significance testing typically used? ( $H_0$ : skeptical reaction to a research hypothesis; ST is used to check if differences or effects seen in the data can be explained simply in terms of sampling variation!)

Slide 14 STAT XL 10, UCLA, Ivo Dinov

## Comments

- How can researchers try to demonstrate that effects or differences seen in their data are real? (Reject the hypothesis that there are no effects)
- How does the alternative hypothesis typically relate to a belief, hunch, or research hypothesis that initiates a study? ( $H_1=H_a$ : specifies the type of departure from the null-hypothesis,  $H_0$  (skeptical reaction), which we are expecting (research hypothesis itself).
- In the Cavendish's mean Earth density data, null hypothesis was  $H_0 : \mu = 5.517$ . We suspected bias, but not bias in any specific direction, hence  $H_a : \mu \neq 5.517$ .

Slide 15 STAT XL 10, UCLA, Ivo Dinov

## Comments

- In the ESP Pratt & Woodruff data, (skeptical reaction) null hypothesis was  $H_0 : \mu = 0.2$  (pure-guessing). We suspected bias, toward success rate being higher than that, hence the (research hypothesis)  $H_a : \mu > 0.2$ .
- Other commonly encountered situations are:
  - $H_0 : \mu_1 - \mu_2 = 0 \rightarrow H_a : \mu_1 - \mu_2 > 0$
  - $H_0 : \mu_{rest} - \mu_{activation} = 0 \rightarrow H_a : \mu_{rest} - \mu_{activation} \neq 0$

Slide 16 STAT XL 10, UCLA, Ivo Dinov

## The t-test

- **Step 1:** Calculate the test-statistic (this tells us how many SD's the estimate is above/below the hypothesized value of the parameter of interest
 
$$t_o = \frac{\hat{\theta} - \theta_0}{SE(\hat{\theta})} = \frac{\text{Estimate} - \text{Hypothesized Value}}{\text{Standard Error}}$$
- **Step 2:** Calculate the P-value from tables or using online resources (e.g., the SOCR, we have online at the class page)
- **Step 3:** Interpret the P-value in context of the data

Slide 17 STAT XL 10, UCLA, Ivo Dinov

## The t-test

Alternative hypothesis	Evidence against $H_0 : \theta > \theta_0$ provided by	P-value
$H_1 : \theta > \theta_0$	$\hat{\theta}$ too much bigger than $\theta_0$ (i.e., $\hat{\theta} - \theta_0$ too large)	$P = \text{pr}(T \geq t_o)$
$H_1 : \theta < \theta_0$	$\hat{\theta}$ too much smaller than $\theta_0$ (i.e., $\hat{\theta} - \theta_0$ too negative)	$P = \text{pr}(T \leq t_o)$
$H_1 : \theta \neq \theta_0$	$\hat{\theta}$ too far from $\theta_0$ (i.e., $ \hat{\theta} - \theta_0 $ too large)	$P = 2 \text{pr}(T \geq  t_o )$

where  $T \sim \text{Student}(df)$

Slide 18 STAT XL 10, UCLA, Ivo Dinov

### Interpretation of the size of the p-value

Approximate size of P-Value	Translation
> 0.12 (12%)	No evidence against $H_0$
0.10 (10%)	Weak evidence against $H_0$
0.05 (5%)	Some evidence against $H_0$
0.01 (1%)	Strong evidence against $H_0$
0.001 (0.1%)	Very Strong evidence against $H_0$

Slide 19 STAT XL 10, UCLA, Ivo Dinov

### Paired Comparisons

- Sometimes we have two data sets, which are not independent, but rather observations matched in pairs.
- When are paired data are significantly different?
- Does the moon size appear different with eyes level and with eyes raised? Does eye position make a difference? Eyes elevated refers to raising the eye from horizontal to zenith position. 10 Subjects are tested under eye-level (control) condition, by physically moving the subject's body from level to zenith position with fixed eye direction - horizontal. Ratios of the Moon size in level and zenith positions, for the two paradigms (physically moving subject's body) are given in Table:

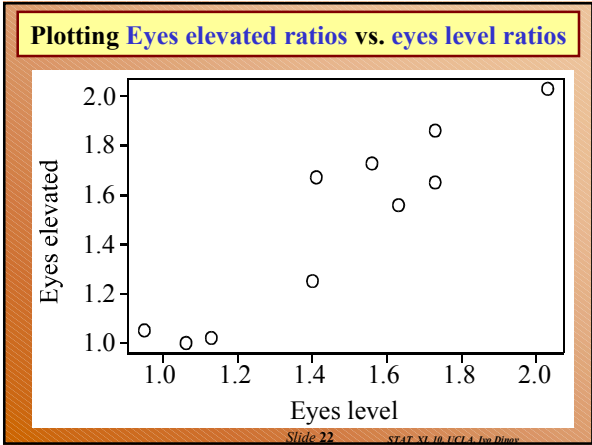
Slide 20 STAT XL 10, UCLA, Ivo Dinov

### Moon illusion Data

Subject	Eyes Elevated	Eyes Level	Difference (Elevated - Level)
1	2.03	2.03	0.00
2	1.65	1.73	-0.08
3	1.00	1.06	-0.06
4	1.25	1.40	-0.15
5	1.05	0.95	0.10
6	1.02	1.13	-0.11
7	1.67	1.41	0.26
8	1.86	1.73	0.13
9	1.56	1.63	-0.07
10	1.73	1.56	0.17

Source: Kaufman and Rock [1962].

Slide 21 STAT XL 10, UCLA, Ivo Dinov



Slide 22 STAT XL 10, UCLA, Ivo Dinov

### Looking for an effect due to elevating eyes

For *paired* data, *analyze the differences*.  $H_0: \mu_{diff} = 0$

Dot plot of differences for the moon illusion data (with a 95% CI for the mean difference).

Variable	N	Mean	StDev	SE Mean	t-stat	P-value
Difference	10	0.0190	0.1371	0.0434	0.44	0.34

Test of  $\mu = 0.0000$  vs  $\mu > 0.0000$   
95% CI ( -0.0791, 0.1171)

Slide 23 STAT XL 10, UCLA, Ivo Dinov

### Comparing two means for independent samples

Suppose we have 2 samples/means/distributions as follows:  $\{\bar{x}_1, N(\mu_1, \sigma_1)\}$  and  $\{\bar{x}_2, N(\mu_2, \sigma_2)\}$ . We've seen before that to make inference about  $\mu_1 - \mu_2$  we can use a T-test for  $H_0: \mu_1 - \mu_2 = 0$  with

$$t_0 = \frac{(\bar{x}_1 - \bar{x}_2) - 0}{SE(\bar{x}_1 - \bar{x}_2)}$$

And  $CI(\mu_1 - \mu_2) = \bar{x}_1 - \bar{x}_2 \pm t \times SE(\bar{x}_1 - \bar{x}_2)$

If the 2 samples are independent we use the SE formula

$$SE = \sqrt{s_1^2/n_1 + s_2^2/n_2} \text{ with } df = \text{Min}(n_1 - 1; n_2 - 1) .$$

Slide 24 STAT XL 10, UCLA, Ivo Dinov

### Means for independent samples – equal or unequal variances?

Pooled T-test is used for samples with assumed equal variances. Under data Normal assumptions and equal variances of  $(\bar{x}_1 - \bar{x}_2 - 0) / SE(\bar{x}_1 - \bar{x}_2)$ , where

$$SE = s_p \sqrt{1/n_1 + 1/n_2}; s_p^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}$$

is exactly Student's *t* distributed with  $df = (n_1 + n_2 - 2)$

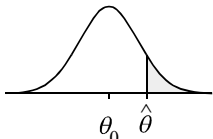
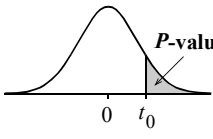
Here  $s_p$  is called the pooled estimate of the variance, since it pools info from the 2 samples to form a combined estimate of the single variance  $\sigma_1^2 = \sigma_2^2 = \sigma^2$ .

Slide 25 STAT.XL.10.UCLA.Ivo.Dimitov

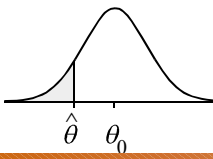
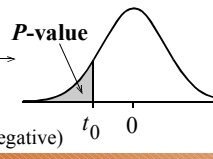
### Comparing two means for independent samples

1. How sensitive is the two-sample *t*-test to non-Normality in the data? (The 2-sample T-tests and CI's are even more robust than the 1-sample tests, against non-Normality, particularly when the shapes of the 2 distributions are similar and  $n_1 = n_2 = n$ , even for small  $n$ , remember  $df = n_1 + n_2 - 2$ .)

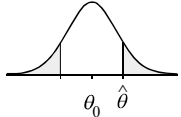
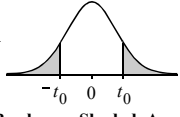
Slide 26 STAT.XL.10.UCLA.Ivo.Dimitov

Alternative Hypothesis	Evidence against $H_0: \theta = \theta_0$ provided by	Pictorial representation of the T-test $H_1: \theta > \theta_0$
$H_1: \theta > \theta_0$	$\hat{\theta}$ too much bigger than $\theta_0$	$t_0 = \frac{\hat{\theta} - \theta_0}{se(\hat{\theta})}$
<b>Back to Hypothesis testing</b> $\hat{\theta}$ -scale $\longrightarrow$ <i>t</i> -scale (# of std errors)		
		

Slide 28 STAT.XL.10.UCLA.Ivo.Dimitov

Alternative Hypothesis	Evidence against $H_0: \theta = \theta_0$ provided by	Pictorial representation of the T-test $H_1: \theta < \theta_0$
$H_1: \theta < \theta_0$	$\hat{\theta}$ too much smaller than $\theta_0$	$t_0 = \frac{\hat{\theta} - \theta_0}{se(\hat{\theta})}$
<b>Back to Hypothesis testing</b> $\hat{\theta}$ -scale $\longrightarrow$ <i>t</i> -scale (# of std errors)		
		

Slide 28 STAT.XL.10.UCLA.Ivo.Dimitov

Alternative Hypothesis	Evidence against $H_0: \theta = \theta_0$ provided by	Pictorial representation of the T-test $H_0: \theta = \theta_0$ $H_1: \theta \neq \theta_0$
$H_1: \theta \neq \theta_0$ (2-sided)	$\hat{\theta}$ too far from $\theta_0$ (either direction)	$t_0 = \frac{\hat{\theta} - \theta_0}{se(\hat{\theta})}$
<b>Back to Hypothesis testing</b> $\hat{\theta}$ -scale $\longrightarrow$ <i>t</i> -scale (# of std errors)		
		

Slide 29 STAT.XL.10.UCLA.Ivo.Dimitov

### P-values from *t*-tests

- The **P-value** is the probability that, if the hypothesis was true, sampling variation would produce an estimate that is further away from the hypothesized value than our data-estimate.
- The **P-value** measures the strength of the evidence against  $H_0$ .
- The **smaller** the **P-value**, the **stronger** the evidence against  $H_0$ .  
(The second and third points are true for significance tests generally, and not just for *t*-tests.)

Slide 31 STAT.XL.10.UCLA.Ivo.Dimitov



### Review

- What does the  $t$ -statistic tell us?  
The T-statistics,  $t_0 = \frac{\theta - \theta_0}{s \alpha(\hat{\theta})}$  tells us (in std. units) if the observed value/estimate is typical/consistent and can be explained by the variation in the sampling distribution.
- When do we use a 2-tailed rather than a 1-tailed test?  
We use two-sided/two-tailed test, unless there is a prior (knowledge available before data was collected) or a strong reason to believe that the result should go in one particular direction ( $\leftarrow \mu \rightarrow$ ).

Slide 32 STAT 10, UCLA, Ivo Dinov

### Review

- What were the 3 types of alternative hypothesis involving the parameter  $\theta$  and the hypothesized value  $\theta_0$ ? Write them down!
- Let's go through and construct our own  $t$ -Test Table.
  - For each alternative, think through what would constitute evidence against the hypothesis and in favor of the alternative.

- Then write down the corresponding  $P$ -values in terms of  $t_0$  and represent these  $P$ -values on hand-drawn curves  
[  $P = \Pr(T > t_0)$ ,  $P = \Pr(T < -t_0)$ ,  $P = 2\Pr(T > |t_0|)$  .]

Slide 33 STAT 10, UCLA, Ivo Dinov

### Review

- What does the  $P$ -value measure? (if  $H_0$  was true, sampling variation alone would produce an estimate farther than the hypothesized value.)
- What do very small  $P$ -values tell us? What do large  $P$ -values tell us? (strength of evidence against  $H_0$ .)
- Pair the phrases: “the  $\uparrow \downarrow$  the  $P$ -value, the  $\uparrow \downarrow$  the evidence for/against the null hypothesis.”
- Do large values of  $t_0$  correspond to large or small  $P$ -values? Why?
- What is the relationship between the Student ( $df$ ) distribution and Normal(0,1) distribution? (identical as  $n \rightarrow \infty$ )

Slide 34 STAT 10, UCLA, Ivo Dinov

### Is a second child gender influenced by the gender of the first child, in families with >1 kid?

First and Second Births by Sex				
		Second Child		Total
		Male	Female	
First Child	Male	3,202	2,776	5,978
	Female	2,620	2,792	5,412
Total		5,822	5,568	11,390

- Research hypothesis needs to be formulated first before collecting/looking/interpreting the data that will be used to address it. Mothers whose 1<sup>st</sup> child is a girl are more likely to have a girl, as a second child, compared to mothers with boys as 1<sup>st</sup> children.
- Data: 20 yrs of birth records of 1 Hospital in Auckland, NZ.

Slide 36 STAT 10, UCLA, Ivo Dinov

### Analysis of the birth-gender data – data summary

Group	Second Child	
	Number of births	Number of girls
1 (Previous child was girl)	5412	2792 (approx. 51.6%)
2 (Previous child was boy)	5978	2776 (approx. 46.4%)

- Let  $p_1$ =true proportion of girls in mothers with girl as first child,  $p_2$ =true proportion of girls in mothers with boy as first child. **Parameter of interest is  $p_1 - p_2$ .**
- $H_0: p_1 - p_2 = 0$  (skeptical reaction).  $H_a: p_1 - p_2 > 0$  (research hypothesis)

Slide 37 STAT 10, UCLA, Ivo Dinov

### Hypothesis testing as decision making

Decision Making		
Decision made	Actual situation	
	$H_0$ is true	$H_0$ is false
Accept $H_0$ as true	OK	Type II error
Reject $H_0$ as false	Type I error	OK

- Sample sizes:  $n_1=5412$ ,  $n_2=5978$ , Sample proportions (estimates)  $\hat{p}_1 = 2792/5412 \approx 0.5159$ ,  $\hat{p}_2 = 2776/5978 \approx 0.4644$ ,
- $H_0: p_1 - p_2 = 0$  (skeptical reaction).  $H_a: p_1 - p_2 > 0$  (research hypothesis)

Slide 38 STAT 10, UCLA, Ivo Dinov

### Analysis of the birth-gender data

- Samples are large enough to use **Normal-approx.**. Since the two proportions come from totally diff. mothers they are **independent** → use formula 8.5.5.a

$$t_0 = \frac{\text{Estimate} - \text{Hypothesized Value}}{SE} = 5.49986 =$$

$$\frac{\hat{p}_1 - \hat{p}_2 - 0}{SE(\hat{p}_1 - \hat{p}_2)} = \frac{\hat{p}_1 - \hat{p}_2}{\sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}} =$$

$$P\text{-value} = \Pr(T \geq t_0) = 1.9 \times 10^{-8}$$

Slide 39 STAT.XI.10.UCLA, Im Dimov

### Analysis of the birth-gender data

- We have strong evidence to reject the  $H_0$ , and hence conclude mothers with first child a girl a **more likely** to have a girl as a second child.

- How much more likely? **A 95% CI:**

CI  $(p_1 - p_2) = [0.033; 0.070]$ . And computed by:

$$\text{estimate} \pm z \times SE = \hat{p}_1 - \hat{p}_2 \pm 1.96 \times SE(\hat{p}_1 - \hat{p}_2) =$$

$$\hat{p}_1 - \hat{p}_2 \pm 1.96 \times \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}} =$$

$$0.0515 \pm 1.96 \times 0.0093677 = [3\%; 7\%]$$

Slide 40 STAT.XI.10.UCLA, Im Dimov

### Review

- Why is the expression “**accept the null hypothesis**” dangerous? (Data can not really provide all the evidence that a hypothesis is true, however, it can provide support that it is false. That’s why better lingo is “**we can’t reject  $H_0$** ”)
- What is meant by the word **non-significant** in many research literatures? (P-value > fixed-level of significance)
- In fixed-level testing, what is a Type I error? What is a Type II error? (Type I, false-positive, reject  $H_0$  as false, when it’s true in reality; Type II, false-negative, accepting  $H_0$  as true, when its truly false)

Slide 44 STAT.XI.10.UCLA, Im Dimov

### Tests and confidence intervals

A **two-sided** test of  $H_0: \theta = \theta_0$  is **significant** at the 5% level **if and only if**  $\theta_0$  lies **outside** a 95% confidence interval for  $\theta$ .

Slide 45 STAT.XI.10.UCLA, Im Dimov

### “Significance”

- **Statistical significance** relates to the strength of the evidence of existence of an effect.
- The **practical significance** of an effect depends on its size – how large is the effect.
- A small P-value provides **evidence that the effect exists** but says **nothing** at all about the **size** of the effect.
- To estimate the **size** of an effect (its practical significance), **compute a confidence interval**.

Slide 47 STAT.XI.10.UCLA, Im Dimov

### “Significance”

- **Statistical significance** relates to the strength of the evidence of existence of an effect, Recall Child-birth example,  $p \sim 2 \times 10^{-8}$ .
- The **practical significance** of an effect depends on its size – how large is the effect. To estimate the **size** of an effect (its practical significance), **compute a confidence interval**. [3%, 7%] **more likely to have a girl as a second child, given the first child is a girl.**

Slide 48 STAT.XI.10.UCLA, Im Dimov

### “Significance” cont.

A non-significant test does not imply that the null hypothesis is true (or that we accept  $H_0$ ).

It simply means we do not have (this data does not provide) the evidence to reject the skeptical reaction,  $H_0$ .

To prevent people from misinterpreting your report: *Never quote a P-value* about the existence of an effect *without* also *providing a confidence interval* estimating the size of the effect.

Slide 49 STAT XI 10, UCLA, Im Dimer

### Review

- What is the relationship between a 95% confidence interval for a parameter  $\theta$  and the results of a two-sided test of  $H_0: \theta = \theta_0$ ? ( $\theta_0$  is inside the 95% CI( $\theta$ ),  $\leftarrow \rightarrow$  P-value for the test is  $> 0.025$ . Conversely, the test is significant, at 5%-level,  $\leftarrow \rightarrow \theta_0$  is outside the 95% CI( $\theta$ ).
- If you read, “research shows that ..... <sup>$\theta$</sup>  is significantly <sup>bigger</sup> than ... $\theta_0$ ..”, what is a likely explanation? (there is evidence that a real effect exists to make the two values different).
- If you read, “research says that .....<sup>drug</sup> makes no difference to .....”, what is a likely explanation? (the data does not have the evidence to reject the skeptical reaction,  $H_0$ ).

Slide 50 STAT XI 10, UCLA, Im Dimer

### Review

- Is a “significant difference” necessarily large or practically important? Why? (No, significant difference indicates the existence of an effect, practical importance depends on the effect-size.)
- What is the difference between statistical significance and practical significance? (stat-significance relates to the strength of the evidence that a real effect exists (e.g., that true difference is not exact,  $y \neq 0$ ); practical significance indicates how important the observed difference is in practice, how large is the effect.)
- What does a P-value tell us about the size of an effect? (P-value says whether the effect is significant, but says nothing about its size.)
- What tool do we use to gauge the size of an effect? (CI(parameter) provides clues to the size of the effect.)

Slide 51 STAT XI 10, UCLA, Im Dimer

### Review

- If we read that a difference between two proportions is *non-significant*, what does this tell us? What does it not tell us? (Do not have evidence proportions are different, based on this data. Doesn't mean accept  $H_0$ ).
- What general strategy can we use to help prevent misconceptions about the meanings of *significance* and *non-significance*? (No, significant difference indicates the existence of an effect, practical importance depends on the effect-size.)
- What is the closest you can get to showing that a hypothesized value is true and how could you go about it? (Suppose,  $H_0: \theta = \theta_0$ , and our test is not-significant. To show  $\theta = \theta_0$  we need to show that all values in the CI( $\theta_0$ ) are essentially equal to  $\theta_0$ , this is a practical subjective matter decision, not a statistical one.)

Slide 52 STAT XI 10, UCLA, Im Dimer

### General ideas of “test statistic” and “p-value”

A *test statistic* is a measure of discrepancy between what we see in data and what we would expect to see if  $H_0$  was true.

The *P-value* is the probability, calculated assuming that the null hypothesis is true, that sampling variation alone would produce data which is more discrepant than our data set.

Slide 53 STAT XI 10, UCLA, Im Dimer

### Course Material Review

- =====Part I=====
- Experiments vs. Observational studies, causality.
- Histograms, dot-plots, stem-and-leaf plot, density curves.
- Numerical summaries of data (5-#-summary)
- The Normal Curve and Normal Approximation
- Percentiles, quartile and linear transformations

Slide 54 STAT XI 10, UCLA, Im Dimer



## Course Material Review – cont.

- Correlation and Regression
- Least squares – best-linear-fit, Linear models
- =====Part II=====
- Probability and proportions (Binomial distribution)
- Confidence Intervals (mean, prop's, & differences)
- Central Limit Theorem
- Hypothesis testing
- Paired vs. Independent samples

Slide 55 STAT XL 10, UCLA, Ivo Dinov

## Chapter 26 – Summary

STAT XL 10, UCLA, Ivo Dinov Slide 56

## Significance Tests vs. Confidence Intervals

- The chief use of **significance testing** is to check whether apparent differences or effects seen in data can be explained away simply in terms of **sampling variation**. The essential **difference between confidence intervals and significance tests** is as follows:
  - **Confidence interval** : A range of possible values for the parameter are determined that are consistent with the data at a specified confidence level.
  - **Significance test** : Only one possible value for the parameter, called the hypothesized value, is tested. We determine the strength of the evidence provided by the data against the proposition that the hypothesized value is the true value.

Slide 57 STAT XL 10, UCLA, Ivo Dinov

## Hypotheses

- The **null hypothesis**, denoted by  $H_0$ , is the (skeptical reaction) hypothesis tested by the statistical test.
- **Principle guiding the formulation of null hypotheses**: We **cannot rule a hypothesized value in**; we can only determine whether there is **enough evidence to rule it out**. **Why is that?**
- **Research (alternative) hypotheses** lay out the conjectures that the research is designed to investigate and, if the researchers hunches prove correct, establish as being true.

Slide 58 STAT XL 10, UCLA, Ivo Dinov

## Hypotheses cont.

- The **null hypothesis** tested is typically a skeptical reaction to the research hypothesis.
- The most commonly tested null hypotheses are of the “it makes no difference” variety.
- Researchers try to demonstrate the existence of real treatment or group differences by showing that the idea that there are no real differences is implausible.
- The **alternative hypothesis**, denoted by  $H_1$ , specifies the type of **departure** from the null hypothesis,  $H_0$ , that we expect to detect.

Slide 59 STAT XL 10, UCLA, Ivo Dinov

## Hypotheses cont.

- The **alternative hypothesis**, typically corresponds to the research hypothesis.
- We use **one-sided alternatives** (using either :  $H_1: \theta > \theta_0$  or  $H_1: \theta < \theta_0$ ) when the research hypothesis specifies the **direction of the effect**, or more generally, when the investigators had good grounds for believing the true value of  $\theta$  was on one particular side of  $\theta_0$  before the study began. Otherwise a **two-sided alternative**,  $H_1: \theta \neq \theta_0$ , is used.

Slide 60 STAT XL 10, UCLA, Ivo Dinov

### P-values

- Differences or effects seen in data that are **easily explainable in terms of sampling variation** **do not provide convincing evidence** that real differences or effects exist.
- The **P-value** is the probability that, if the hypothesis was true, sampling variation would produce an estimate that is further away from the hypothesized value than the estimate we got from our data.
- The **P-value** **measures the strength of the evidence against  $H_0$** .

Slide 61 STAT.XL.10.UCLA.Ivo.Dimitrov

### P-values cont.

- The *smaller* the **P-value**, the stronger the evidence against  $H_0$ .
- A large **P-value** provides no evidence against the null hypothesis.
- A large **P-value** does *not* imply that the null hypothesis is true.
- A small **P-value** provides evidence that the effect exists but says *nothing* at all about the *size* of the effect.
- To estimate the **size** of an effect, *compute a confidence interval*.

Slide 62 STAT.XL.10.UCLA.Ivo.Dimitrov

### P-values cont.

- Never quote a **P-value** about the existence of an effect without also providing a confidence interval estimating the size of the effect.
- Suggestions for **verbal translation of P-values** are given in Table 9.3.2.
- **Computation of P-values** : Computation of P-values for situations in which the sampling distribution of  $(\hat{\theta} - \theta_0) / se(\hat{\theta})$ , is well **approximated by a Student(df) distribution or a Normal(0,1)** distribution is laid out in Table 9.3.1.
- The **t-test** statistic tells us how many standard errors the estimate is from the hypothesized value.

Slide 63 STAT.XL.10.UCLA.Ivo.Dimitrov

### P-values

- Examples given in this chapter concerned means and differences between means, proportions and differences between proportions.
- In general, a test statistic is a measure of discrepancy between what we see in the data and what we would have expected to see if  $H_0$  was true.

Slide 64 STAT.XL.10.UCLA.Ivo.Dimitrov

### Significance

- If, whenever we obtain a **P-value** less than or equal to 5%, we make a decision to reject the null hypothesis, this procedure is called **testing at the 5% level of significance**.
  - The significance level of such a test is 5%.
- If the **P-value**  $\leq \alpha$ , the effect is said to be significant at the  $\alpha$ -level.
- If you always test at the 5% level, you will reject one true null hypothesis in 20 over the long run.

Slide 65 STAT.XL.10.UCLA.Ivo.Dimitrov

### Significance cont.

- A two-sided test of  $H_0 : \theta = \theta_0$  is significant at the 5% level if and only if  $\theta_0$  lies outside a 95% confidence interval for  $\theta$ .
- In reports on research, the word “**significant**” used alone often means “**significant at the 5% level**” (i.e. **P-value**  $\leq 0.05$ ). “**Non-significant**”, “**does not differ significantly**” and even “**is no different**” often mean **P-value**  $> 0.05$ .
- A non-significant result does not imply that  $H_0$  is true.

Slide 66 STAT.XL.10.UCLA.Ivo.Dimitrov

### Significance cont.

- A Type I error (false-positive) is made when one concludes that a true null hypothesis is false.
- The significance level is the probability of making a Type I error.
- *Statistical significance* relates to having evidence of the *existence* of an effect.
- The *practical significance* of an effect depends on its *size*.