**Chapter 3: A Closer Look at Assumptions [of the *t* tests]
Chapter 4: Alternatives to the t tools [2 full classes]**

2/18/09 W

**Slide 1    Chapter 3: A Closer Look at Assumptions [of the t tests]**

Chapter 4: Alternatives to the t tools [2 full classes]

NOTES:

---

**HW 4 due Fri 2/20/09 Noon**

Submit as Myname-HW4.doc (or *.rtf)

- Finish Ch 3 for Weds' class
  ‣ Chapter 3: A closer look at assumptions
  ‣ Read
    ▪ Hayek & Buzas (1997, on sampling)
    ▪ Hurlbert (1984) on Pseudoreplication
    ▪ **Post one comment and one reply to issues raised in Hayek & Buzas or Hurlbert (1984)**
- Chapter 3 problem due Weds 2/18
  ‣ **3.28 Pollen removal**
- Read all of chapter 4: Wilcoxon rank sum, signed rank tests, Fisher's sign test, Welch's unequal variance t test

**Slide 2  HW 4 due Fri 2/20/09 Noon**

NOTES:

---

**HW 5 due Weds 2/25/09 9:50**

Submit as Myname-HW5.doc (or *.rtf)

- Finish Chapter 4 **Wilcoxon rank sum, signed rank tests, Fisher's sign test, Welch's unequal variance t test**
- Comment on Chapter 4 conceptual problems in Blackboard Vista4
- Computation Problem 5
  ‣ Problem 4.31 Effect of group therapy on breast cancer patients.

**Slide 3  HW 5 due Weds 2/25/09 9:50**

NOTES:

**HW 6 due Monday 3/1/09 9:50**

Submit as Myname-HW5.doc (or *.rtf)

- Read Chapter 5 **Comparisons among several samples**
- Comment on Chapter 5 conceptual problems in Blackboard Vista4
- Computation Problem 6
  - Problem 4.30 Sunlight protection factor

| Slide 4  HW 6 due Monday 3/1/09 9:50 |
| --- |
|  |
| NOTES: |
|  |
|  |
|  |
|  |

**Chapter 3: A closer look at assumptions**

| Slide 5  Chapter 3: A closer look at assumptions |
| --- |
|  |
| NOTES: |
|  |
|  |
|  |
|  |
|  |

**Case study 3.1**

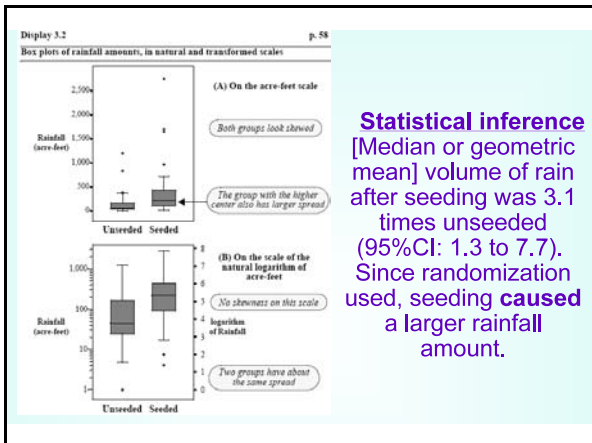Cloud seeding to increase rainfall — A randomized experiment

- 52-day experiment
- Random selection each day to seed or not to seed a cloud; pilot 'blind' to treatment
- Rainfall measured
- Data highly skewed

Display 3.1

Rainfall (acre-feet) for days with and without cloud seeding

Rainfall from unseeded days (n = 26)

| 1202.6 | 830.1 | 372.4 | 345.5 | 321.2 | 244.3 | 163.0 | 147.8 | 95.0 |
| 87.0 | 81.2 | 68.5 | 47.3 | 41.1 | 36.6 | 29.0 | 28.6 | 26.3 |
| 26.1 | 24.4 | 21.7 | 17.3 | 11.5 | 4.9 | 4.9 | 1.0 | |

Rainfall from seeded days (n = 26)

| 2745.6 | 1697.8 | 1656.0 | 978.0 | 703.4 | 489.1 | 430.0 | 334.1 | 302.8 |
| 274.7 | 274.7 | 255.0 | 242.5 | 200.7 | 198.6 | 129.6 | 119.0 | 118.3 |
| 115.3 | 92.4 | 40.6 | 32.7 | 31.4 | 17.5 | 7.7 | 4.1 | |

| Slide 6  Case study 3.1 |
| --- |
|  |
| NOTES: |
|  |
|  |
|  |
|  |

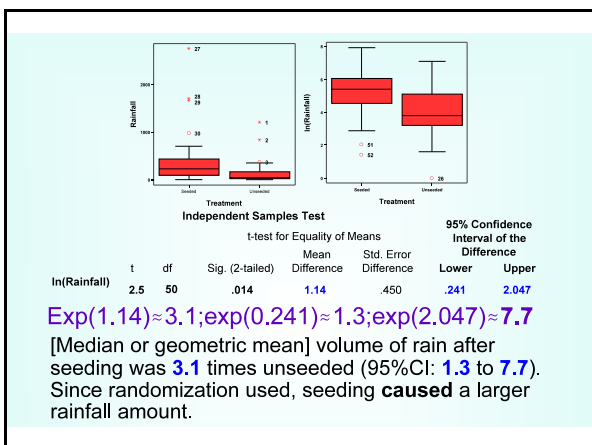| | |
|---|---|
| **Display 3.2** p. 58 — Box plots of rainfall amounts, in natural and transformed scales. **Statistical inference** [Median or geometric mean] volume of rain after seeding was 3.1 times unseeded (95%CI: 1.3 to 7.7). Since randomization used, seeding **caused** a larger rainfall amount. | **Slide 7**<br><br>NOTES: |
| Exp(1.14)≈3.1;exp(0.241)≈1.3;exp(2.047)≈7.7 [Median or geometric mean] volume of rain after seeding was **3.1** times unseeded (95%CI: **1.3** to **7.7**). Since randomization used, seeding **caused** a larger rainfall amount. | **Slide 8**<br><br>NOTES: |
| **Randomization doesn't solve problems with unequal variance** | **Slide 9  Randomization doesn't solve problems with unequal variance**<br><br>NOTES: |

## Case 3.2: Dioxin study

**Differences between veteran dioxin concentrations could be due to chance**

- **646 Veterans who served in Viet Nam during 1967 & 1968 in areas treated with Agent Orange**
- **97 other veterans served between 1965-1971 in US or Germany**
- **Serum dioxin levels measured**
- **Statistical Summary:**
  - No evidence that the mean dioxin levels differ (1-sided p value=0.4)
  - Extrapolation speculative; dioxin-affected vets may not have participated in the survey

*EEOS611*

**Slide 10  Case 3.2: Dioxin study**

NOTES:

---

## Robustness of the two-sample *t* tools

**Slide 11  Robustness of the two-sample t tools**

NOTES:

---

## Assumptions of *t* test

- Two major assumptions
  - Both samples are independent samples from normally distributed populations
  - Both samples have identical standard deviations
- The *t* tests are usually robust to modest violations of the assumptions
  - These assumptions are never strictly met, but the *t* test is remarkably robust to violations of the assumptions
  - Robust means the conclusions from test — *e.g., p* values, confidence limits — are valid even when the assumptions aren't strictly met, especially if sample sizes nearly equal
  - Transformations of the data are often used

*EEOS611*

**Slide 12  Assumptions of t test**

NOTES:

## Violations of assumptions that matter

- With equal sample sizes, the *t*-test is affected moderately by long-tailedness (leptokurtic or peaked distribution) and very little by skewness (the symmetry of the distribution)
  ‣ Kurtosis: peakedness, platykurtic (flat distribution), leptokurtic (peaked)
  ‣ **Skewness: symmetry**
- If the two populations have the same standard deviations and approximately the same shape, with unequal sample size, the *t* tests are affected moderately by long tailedness (leptokurtic) and **substantially by skewness**
- If the **skewness** differs considerably, the **tools can be misleading with small and moderate sample sizes**
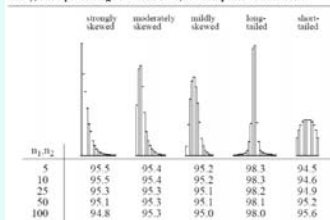
*EEOS611*

**Slide 13  Violations of assumptions that matter**

NOTES:

## Monte Carlo simulations of violations on p values

Display 3.4    Which violations of assumptions really matter?

Percentage of 95% confidence intervals that are successful when the two populations are non-normal (but same shape and SD, and equal sample sizes); each percentage is based on 1,000 computer simulations

| n1,n2 | strongly skewed | moderately skewed | mildly skewed | long-tailed | short-tailed |
|---|---|---|---|---|---|
| 5 | 95.5 | 95.4 | 95.2 | 98.3 | 94.5 |
| 10 | 95.5 | 95.4 | 95.2 | 98.3 | 94.6 |
| 25 | 95.3 | 95.3 | 95.1 | 98.2 | 94.9 |
| 50 | 95.1 | 95.3 | 95.1 | 98.1 | 95.2 |
| 100 | 94.8 | 95.3 | 95.0 | 98.0 | 95.6 |

When 2 populations are the same shape with equal n, the results of the *t* test affected moderately (conservative for leptokurtic [peaked] distributions)

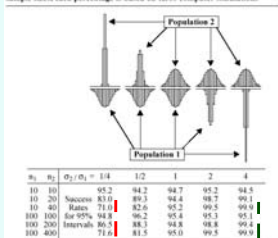**Slide 14  Monte Carlo simulations of violations on p values**

NOTES:

## Different standard deviations & sample sizes

Robust if sample sizes the same, nonconservative if ...ger s.d.

Display 3.5    Percentage of successful 95% confidence intervals when the two populations have different standard deviations (but are normal) with possibly different sample sizes; each percentage is based on 1,000 computer simulations

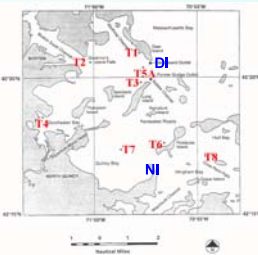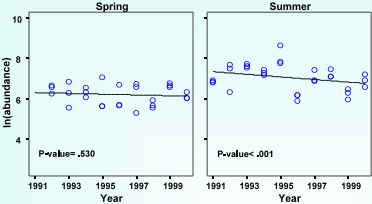| n1 | n2 | σ2/σ1 = 1/4 | 1/2 | 1 | 2 | 4 |
|---|---|---|---|---|---|---|
| 10 | 10 | 95.2 | 94.2 | 94.7 | 95.2 | 94.5 |
| 10 | 20 | Success 83.0 | 89.3 | 94.4 | 98.7 | 99.1 |
| 10 | 40 | Rates 71.0 | 82.6 | 95.2 | 99.5 | 99.9 |
| 100 | 100 | for 95% 94.8 | 96.2 | 95.4 | 95.3 | 95.1 |
| 100 | 200 | Intervals 86.5 | 88.3 | 94.8 | 98.8 | 99.4 |
| 100 | 400 | 71.6 | 81.5 | 95.0 | 99.5 | 99.9 |

P (Type I error) >> stated value (e.g., 0.05) if sd of smaller group larger than larger group

P (Type I error) << stated value (e.g., 0.05) if sd of smaller group smaller than larger group

*EEOS611*

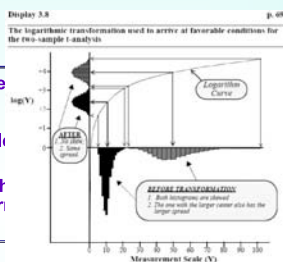**Slide 15  Different standard deviations & sample sizes**

NOTES:

## Departures from independence

**Cluster, serial & spatial effects can be serious and are more difficult to account for in statistical analysis**

- Cluster effects
  - Mice from litters
  - Copepods from net hauls
  - Technician-to-technician variability in sample analysis
- Serial effects (an explicit time or space term)
  - Temporal autocorrelation
- Spatial effects: positive autocorrelation (Ellen Douglas's research on flood frequency, Chen & Ferguson on MCAS scores) **EEOS611**
- Check residuals (observed-expected) for spatial or temporal pattern
- Inferences based on Student's *t* tests can be very misleading or wrong if there are spatial

| **Slide 16  Departures from independence** |
|---|
| |
| NOTES: |
| |
| |
| |
| |
| |

## MWRA Benthic Sampling Stations

**8 Stations, May & Aug**
**3 replicate 0.043-m² Ted Young grabs; 300-µm sieves**

- T1: Deer Island
- T2: Governor's Island Flats
- T3: Long Island
- T4: Savin Hill Cove
- T5A: Presidents Road
- T6: Peddocks Island
- T7: Quincy Bay
- T8: Hingham/Hull Bay
- NI: Nut Island
- DI: Deer Island

| **Slide 17  MWRA Benthic Sampling Stations** |
|---|
| |
| NOTES: |
| |
| |
| |
| |
| |

The residuals after fitting the regression should be identically independently normally distributed, but they are not. The analyst can not ignore these effects and perform a regression as if these residuals are independent

Fig. 40. The change in ln(abundance) at Quincy Bay(T07). There is a significant lack of fit in both the spring and summer data to linear regression. The p-values reported are from the One-way ANOVA test. Banik 2003 UMB M.Sc.

Problems with serial autocorrelation (confounded with spatial effects) create a problem called 'lack of fit' in OLS and regression

Solution: ANOVA test for linear trends

| **Slide 18** |
|---|
| |
| NOTES: |
| |
| |
| |
| |
| |

| | |
|---|---|
| **Log transform of rainfall** | **Slide 19  Log transform of rainfall** |
| See Case 3.1 movie | |
| Display 3.8                                    p. 69 | NOTES: |
| The logarithmic transformation used to arrive at favorable conditions for the two-sample t-analysis | |
| ●The antilogarithm of the me... of the log values, the geometric mean, is the median on the original scale of measurement | |
| ●Calculate the 95% CI's on th... log scale and back transfor... they will be asymmetric | |
| *EEOS611* | |

| | |
|---|---|
| **3.5.3 Transformations** | **Slide 20  3.5.3 Transformations** |
| ● Log (x+1) transform | |
| ▸ Most biological data, but not usually diversity | NOTES: |
| ▸ Needed when there is a multiplicative process in action: growth, bank account interest | |
| ▸ Marine pollutants: polynuclear aromatic hydrocarbons, fecal coliform bacteria, but not usually metals | |
| ▸ Calculate the mean and 95% CI and then back-transform. For symmetric data, the mean of the log-transformed data ≈median. Label as the geometric mean | |
| ● Many other transforms | |
| ▸ Arcsin ($\sqrt{Y}$) for frequency data ranging between 0 and 1 (but the logit transform may be better) | |
| ▸ % silt clay, but the data must be on the interval 0 to1 | |
| ▸ Logit transform: log [Y/(1-Y)] | |
| ▸ Square roots for counts, reciprocal for waiting times, logit transforms for proportions between 0 and 1 (log (P/(1-P)) | |
| ▸ "... it is recommended here that a trial-and error approach, with graphical analysis, be used instead." | |

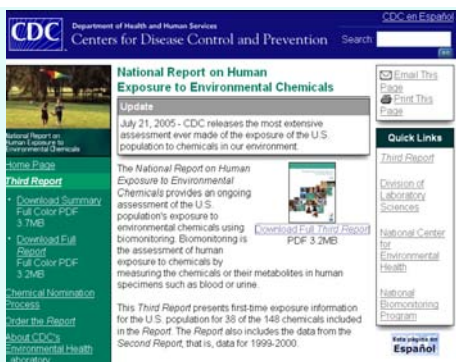| | |
|---|---|
| Display 3.8                                    p. 68 | **Slide 21** |
| Two-sample t-analysis and statement of conclusions after logarithmic transformation — cloud seeding example | |
| ① Transform the data | NOTES: |
| Do the test and calculate the 95% CI on transformed data and then back transform the effect size and confidence limits. Report as ratio of geometric means (Sleuth: ratio of medians). | |

| | |
|---|---|
| <br>http://www.cdc.gov/exposurereport/ | **Slide 22**<br><br>NOTES:<br><br><br><br><br> |
| <br>http://www.cdc.gov/exposurereport/3rd/pdf/thirdreport.pdf | **Slide 23**<br><br>NOTES:<br><br><br><br><br> |
|  | **Slide 24**<br><br>NOTES:<br><br><br><br><br> |

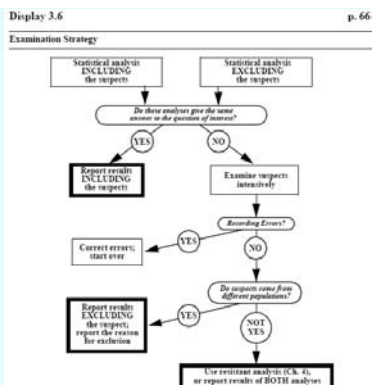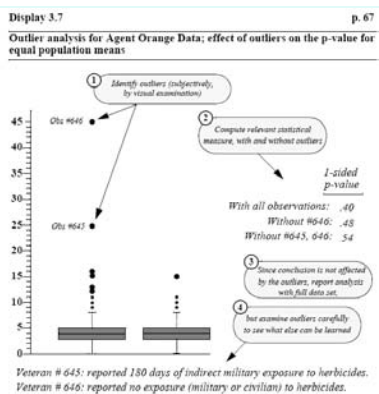| | |
|---|---|
| **Outliers and resistant procedures**<br><br>● A procedure is resistant if it doesn't change much when a small part of the data changes, perhaps drastically.<br>● *t* tools are based on averages and are strongly affected by outliers<br>‣ Chapter 4 introduces tests based on ranks, which protect against outliers (but not against unequal variance)<br>● Practical strategies<br>‣ Do side-by-side box plots to analyze departures from assumptions<br>  ■ Check for patterns in residuals with box plots<br>‣ Consider & test for serial spatial and cluster effects<br>  ■ Analyze spatial patterns in the residuals, use more sophisticated tools<br>  ■ Legendre & Legendre: if pos. Spatial autocorrelation, decrease the p value. Test for differences at the 0.001 level instead of the 0.05 level | **Slide 25  Outliers and resistant procedures**<br><br>NOTES: |

| | |
|---|---|
| <br>Display 3.6     p. 66<br>**Examination Strategy** | **Slide 26**<br><br>NOTES: |

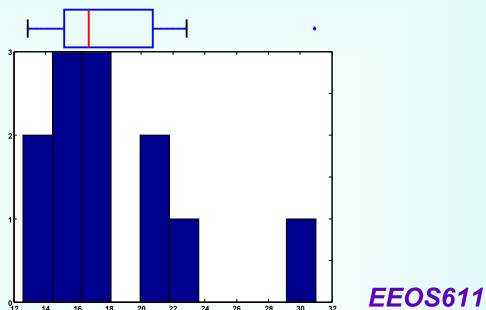| | |
|---|---|
| <br>Display 3.7     p. 67<br>Outlier analysis for Agent Orange Data; effect of outliers on the p-value for equal population means<br><br>Report results, with and without outliers<br><br>Veteran # 645: reported 180 days of indirect military exposure to herbicides.<br>Veteran # 646: reported no exposure (military or civilian) to herbicides. | **Slide 27  Identifying outliers with boxplots**<br><br>NOTES: |

## Practical strategies for outliers

**Be wary of outlier deletion!**

- Outlier strategy
  - ‣ Run analysis with and without outliers
  - ‣ Throw-out outliers only if there is very compelling evidence to do so, and document this data paring or culling
- Note that outlier removal has created tremendous problems:
  - ‣ POC flux to the deep sea
  - ‣ The ozone hole
  - ‣ Mendel's data: 1:2:1 ratios and the chi-square test; documented by Fisher
  - ‣ Milliken's study of the charge of the electron

**Slide 28  Practical strategies for outliers**

NOTES:

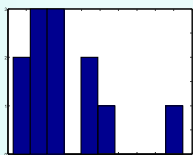## Is the datum at 30 an outlier?



*EEOS611*

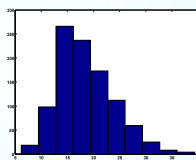**Slide 29  Is the datum at 30 an outlier?**

NOTES:

## Lognormal distribution

**Mean 18.6, Standard deviation 5, n=1000**

I generated the data from a parametric lognormal distribution with μ=18.6 and σ = 5.
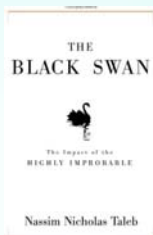


*EEOS611*

**Slide 30  Lognormal distribution**

NOTES:

## Outliers & Black Swans

- Taleb argues that many events are characterized by extreme events
- Mandelbrotian grey swans are events that can be partially characterized through transformations (earthquakes, etc.) But these require data
- Luddian falacy is the belief that all events can be characterized by probabilistic models

THE
BLACK SWAN

The Impact of the
HIGHLY IMPROBABLE

Nassim Nicholas Taleb

*EEOS611*

---

**Slide 31  Outliers & Black Swans**

NOTES:

---

## Chapter 4: Alternatives to the *t* tools

Note: Sleuth has MANY errors and omissions!

- Permutation tests [Not a solution to unequal variance]
- Wilcoxon's Rank Sum Test (same probability model as Mann-Whitney U test)
  ‣ Ties corrections not in sleuth & exact p values
- Repeated Measures Tests based on ranks: Wilcoxon Sign Rank & Fisher's Sign Test
- Parametric vs. Nonparametrics
  ‣ Power efficiency not an issue, ties not that much of an issue
  ‣ Dealing with covariates & estimating effect sizes can be an issue
    ▪ Hodges-Lehman estimators
- Unequal variance (Welch's) t test: some theoretical and practical problems
- Supplemental material
  ‣ Two-sample binomial test (covered in Sleuth Ch 19)
  ‣ The Fligner-Policello test, a rank-based test for samples with unequal variance

---

**Slide 32  Chapter 4: Alternatives to the t tools**

NOTES:

---

## Case 4.1: Space Shuttle O-Ring Failures

See Case 4.1 Movie, solved in Matlab™ & SPSS

Display 4.1

Numbers of O-ring incidents on 24 space shuttle flights prior to the Challenger disaster

*Two problems: unequal variance & ties*

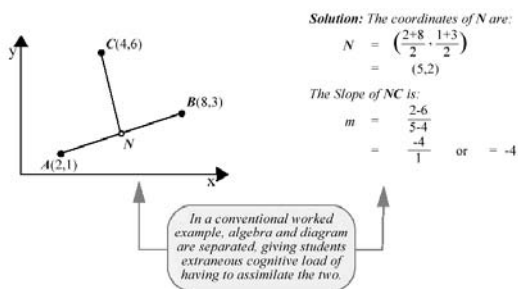| Launch Temperature | Number of O-Ring Incidents |
|---|---|
| Below 65° F | 1 1 1 3 |
| Above 65° F | 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 1 2 |

**Summary of Statistical Findings**
There is strong evidence that the number of O-ring incidents was associated with launch termperature in these 24 launches ...
p- value = 0.009 from a **permutation test** on the *t* statistic

---

**Slide 33  Case 4.1: Space Shuttle O-Ring Failures**

NOTES:

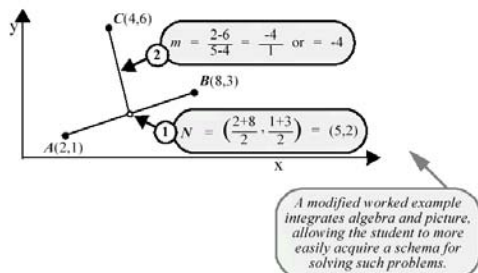| | Slide 34  Case 4.2: Cognitive Load |
|---|---|
| **Display 4.2** Cognitive load experiment: conventional method of instruction (for finding the slope of the line that connects C to the midpoint between A and B)  | NOTES: |

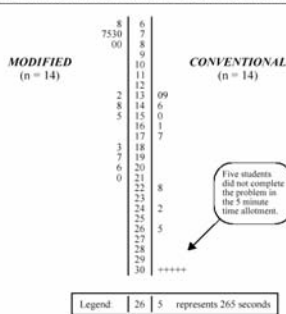| | Slide 35 |
|---|---|
| **Display 4.3** Cognitive load experiment: modified method of instruction (for finding the slope of the line that connects C to the midpoint between A and B)  | NOTES: |

| | Slide 36 |
|---|---|
| **Display 4.4** Numbers of seconds to solution of a problem in coordinate geometry, for students instructed with conventional and modified materials  *Problem: censored (truncated) data* **Statistical Summary** There was convincing evidence that a student could solve the problem more quickly if taught with the modified method (1-sided p-value = 0.003 from the rank-sum test).  The modified method shortened solution times by an estimated 152 s (95% CI 58 to 159 s) | NOTES: |

| | |
|---|---|
|  | **Slide 37  Wilcoxon's rank-sum test** |
| | |
| | NOTES: |
| | |
| | |
| | |
| | |
| | |

Slide content:

**Wilcoxon's rank-sum test**

**Analogous to 2-sample t test**
**Same test as the Mann- Whitney U test**

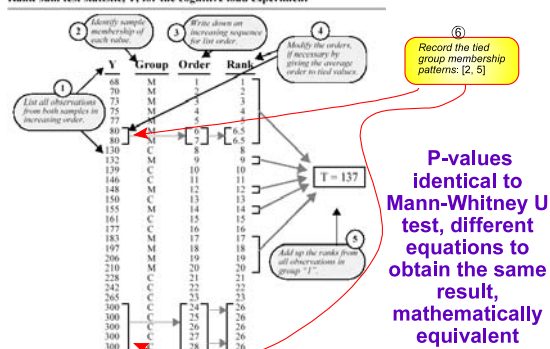Frank Wilcoxon of American Cyanamide [See Salsburg, 2001, The Lady Tasting Tea]

| | |
|---|---|
|  | **Slide 38  Wilcoxon Rank Sum Statistic** |
| | |
| | NOTES: |
| | |
| | |
| | |
| | |
| | |

Slide content (Display 4.5):

Rank-sum test statistic, T, for the cognitive load experiment

Record the tied group membership patterns: [2, 5]

$T = 137$

**P-values identical to Mann-Whitney U test, different equations to obtain the same result, mathematically equivalent**

| | |
|---|---|
|  | **Slide 39** |
| | |
| | NOTES: |
| | |
| | |
| | |
| | |
| | |

Slide content (Display 4.6):

Facts about the randomization (or sampling) distribution of the rank-sum statistic—the sum of ranks in group I—when there is no group difference

**PERMUTATION DISTRIBUTION OF THE RANK-SUM (T)**

③ SHAPE
The shape of the sampling distribution will be approximately normal if the sample sizes are large (and not too many ties)

Correct the standard deviation SD(T), based on the pattern of ties [2, 5]

① CENTER
$Mean(T) = n_1 \overline{R}$

② SPREAD
$SD(T) = s_R \sqrt{\dfrac{n_1 n_2}{(n_1 + n_2)}}$

where $\overline{R}$ and $s_R$ are the average and the sample standard deviation for the combined set of ranks (e.g. the 4th column of Display )

**Note: Sleuth doesn't include the ties correction for the variance of Wilcoxon T**

| | |
|---|---|
| **Display 4.7**<br><br>Finding the p-value with the normal approximation to the permutation distribution of the rank-sum statistic; calculations for the cognitive load data continued from .<br><br>① Calculate the average and sample standard deviation of the ranks from the combined sample (column 4 of Display )<br><br>$\bar{R} = 14.5$    $s_R = 8.202$<br><br>③ **Correct the standard deviation SD(T), based on the pattern of ties [2, 5]**<br><br>② Compute the theoretical "null hypothesis" mean and standard deviation of T, using the formulas in<br><br>$\text{Mean}(T) = (14)(14.5) = 203;$    $SD(T) = 8.202\sqrt{\dfrac{14 \times 14}{(14+14)}} = 21.70$<br><br>③ Determine the Z-statistic.  →  $Z = \dfrac{(137 - 203)}{21.70} = -3.04$<br><br>④ Find the p-value from a standard normal table  →  one-sided p-value = .00118 | **Slide 40**<br><br>NOTES: |

| | |
|---|---|
| **Ties Correction for Rank Sum**<br><br>From Hollander & Wolfe (1999); ties reduce Var(W) and produce a more powerful & accurate test<br><br>**Ties**<br>If there are ties, give tied observations the average of the ranks for which those observations are competing. After computing W using average ranks, use procedures (4.4), (4.5) or (4.6) and refer the value of W to Table A.6. Now, however, the test is approximate rather than exact. (To get an exact test, even in the tied case, see Comment 5.)<br>When applying the large-sample approximation, the following modification should be made. When there are ties, the null mean of W is unaffected, but the null variance is reduced to<br><br>$var_0(W) = \dfrac{mn}{12}\left[m + n + 1 - \dfrac{\sum_{j=1}^{g}(t_j - 1)t_j(t_j + 1)}{(m+n)(m+n-1)}\right],$    (4.13)<br><br>or, equivalently,<br><br>$var_0(W) = \dfrac{mn(N+1)}{12} - \left\{\dfrac{mn}{12N(N-1)} \cdot \sum_{j=1}^{g}(t_j - 1)t_j(t_j + 1)\right\}.$    (4.14)<br><br>To apply the large-sample approximation when ties are present, compute W using average ranks, and compute<br><br>$W^* = \dfrac{W - [m(m + n + 1)/2]}{[var_0(W)]^{1/2}}$<br><br>where $var_0(W)$ is given by display (4.13). With this modified value of $W^*$, approximations (4.10), (4.11) and (4.12) can be applied.<br><br>**g=tied groups**<br>**$t_j$ = items in each tied group**<br><br>*EEOS611* | **Slide 41  Ties Correction for Rank Sum**<br><br>NOTES: |

| | |
|---|---|
| **Mann-Whitney U test**<br><br>Wilcoxon (1945), Mann & Whitney (1947)<br><br>● The statistic U can be computed as follows (Hollander & Wolfe 1999)<br>  ▸ For the two groups $X_i$ and $Y_j$ with m & n cases, consider each of the m×n pairs<br>  ▸ For each pair of values $X_i$ and $Y_j$, observe which is smaller.<br>  ▸ If the $X_i$ value is smaller, <span style="color:red">score a 1 for that pair.</span> If the $Y_j$ value is smaller, score a 0 for that pair.<br>  ▸ Mann & Whitney showed that in the case of no ties:<br>    ■ T=U+[n(n+1)/2], where T is the sum of ranks from the Wilcoxon rank sum test<br>    ■ Thus, the tests are exactly equivalent<br>  ▸ When, $X_i$ and $Y_j$ are tied, <span style="color:red">score ½</span><br><br>*EEOS611* | **Slide 42  Mann-Whitney U test**<br><br>NOTES: |

| | |
|---|---|
|  | <br><br>NOTES: |
|  | **Slide 44**<br><br>NOTES: |
|  | <br><br>NOTES: |

**SPSS's Mann-Whitney U test**

C:\program files\spss\help\algorithms\npar_tests.pdf

**Mann-Whitney *U* Test**

**Calculation of Sums of Ranks**

The combined data from both groups are sorted and ranks assigned to all cases, with average rank being used in the case of ties. The sum of ranks for each of the groups ($S_1$ and $S_2$) is calculated, as well as, for tied observations, $T_i = \frac{t^3 - t}{12}$, where $t$ is the number of observations tied for rank $i$.

The average rank for each group is

$$\bar{S}_i = S_i / n_i$$

where $n_i$ is the sample size in group $i$.

EEOS611

**Test Statistic and Significance Level**

The $U$ statistic for group 1 is

$$U = n_1 n_2 + \frac{n_1(n_1+1)}{2} - S_1$$

- If $U > n_1 n_2/2$, the statistic used is

$$U' = n_1 n_2 - U$$

- If $n_1 n_2 \le 400$ and $n_1 n_2/2 + \min(n_1, n_2) \le 220$ the exact significance level based on an algorithm of Dineen and Blakesley (1973).

- The test statistic corrected for ties is

$$Z = \frac{(U - n_1 n_2/2)}{\sqrt{\frac{n_1 n_2}{N(N-1)}\left(\frac{N^3 - N}{12} - \sum_i T_i\right)}}$$

$T_i$ = number of items in each tied group

which is distributed approximately as a standard normal. A two-tailed significance level is printed.

**Exact p values tabulated (no ties)**

**If no ties, exact P values tabulated**

- CRC Handbook of Tables for Probability and Statistics
  - ‣ Tabulated values of Mann-Whitney U statistic
  - ‣ Can readily convert from sum of ranks of smaller group to U statistics
- Or, Hollander & Wolf's (1999) tabulated values of the Wilcoxon T statistic
- Note: the p values for tabulated exact tests are not appropriate if there are any tied ranks; but an exact p value can be calculated using all combinations of data (Gallagher provides a Matlab m.file implementing Hollander & Wolfe algorithm)

## Slide 46  Other alternatives for two independent samples

**Other alternatives for two independent samples**

4.3.1 Permutation tests

NOTES:

## Slide 47  Ties and Case 4.1

**Ties and Case 4.1**

**Sleuth argues that rank tests inappropriate because of ties. The real problem is unequal variance (Behrens**

Display 4.1

Numbers of O-ring incidents on 24 space shuttle flights prior to the Challenger disaster

| Launch Temperature | Number of O-Ring Incidents |
|---|---|
| Below 65° F | 1 1 1 3 |
| Above 65° F | 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 1 2 |

Sleuth argues that tied values pose problems for the Wilcoxon rank sum test, **but** Siegel argues that the Wilcoxon rank sum test is robust in the presence of ties. The test is only approximate with ties, and the normal approximation is conservative. There is an exact test with ties (computer intensive but Gallagher has programmed).

NOTES:

## Slide 48

Display 4.1

Numbers of O-ring incidents on 24 space shuttle flights prior to the Challenger disaster

| Launch Temperature | Number of O-Ring Incidents |
|---|---|
| Below 65° F | 1 1 1 3 |
| Above 65° F | 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 1 2 |

Independent Samples Test

| | t | df | Sig. (2-tailed) | 95% Confidence Interval of the Difference Lower | Upper |
|---|---|---|---|---|---|
| Equal variances assumed | 3.9 | 22 | .00079 | .61 | 2.0 |
| Equal variances not assumed | 2.5 | 3.34 | .07690 | -.25 | 2.8 |

Test Statistics[b]

| | Incident |
|---|---|
| Mann-Whitney U | 6.000 |
| Wilcoxon W | 216.000 |
| Z | -3.301 |
| Asymp. Sig. (2-tailed) | .000963 |
| Exact Sig. [2*(1-tailed Sig.)] | .005082[a] |

a. Not corrected for ties.

b. Grouping Variable: Temperature

**Note the 100-fold difference in p values**

NOTES:

**Wilcoxon rank sum test assumptions**

Nonparametric tests distribution-free, NOT assumption free

- The observations of $X_1, ..., X_m$ are a random sample from population 1, independent and identically distributed. The observations of $Y_1, ..., Y_m$ are a random sample from population 2, independent and identically distributed
- The X's and Y's are mutually independent
- Populations 1 and 2 are continuous [i.e., no ties]
- **"Robustness of level: The significance level of the rank sum test is not preserved if the two populations differ in dispersion or shape. This is also the case for the normal theory 2-sample t test."** Hollander & Wolfe, p. 120

| **Slide 49  Wilcoxon rank sum test assumptions** |
| --- |
| |
| NOTES: |
| |
| |
| |
| |
| |

---

Display 4.10

A summary of the t-statistics calculated from all 10,626 rearrangements of the O-ring data into a "Low" group of size 4 and a "High" group of size 20

| Number of re-arrangements with identical t-statistics | t-statistic |
| --- | --- |
| 2,380 | -1.188 |
| 3,400 | -0.463 |
| 2,040 | 0.231 |
| 1530 | 0.939 |
| 855 | 1.716 |
| 316 | 2.643 |
| 95 | 3.888 |
| 10 | 5.952 |

Total number of rearrangements into two groups of size 4 and 20:
10,626

Number of rearrangements with t-statistics greater than or equal to 3.888:
105

1-sided p-value from a permutation test of the t-statistic:
105/10626 = .00988

This approach is invalid.  The underlying t test assumes equal variances, and that problem is not corrected by using permutations.  See Manly

| **Slide 50** |
| --- |
| |
| NOTES: |
| |
| |
| |
| |
| |

---

**nCr=24 Choose 4=24!/((24-4)!*4!)**

Display 4.10

A summary of the t-statistics calculated from all 10,626 rearrangements of the O-ring data into a "Low" group of size 4 and a "High" group of size 20

| Number of re-arrangements with identical t-statistics | t-statistic |
| --- | --- |
| 2,380 | -1.188 |
| 3,400 | -0.463 |
| 2,040 | 0.231 |
| 1530 | 0.939 |
| 855 | 1.716 |
| 316 | 2.643 |
| 95 | 3.888 |
| 10 | 5.952 |

Total number of rearrangements into two groups of size 4 and 20:
10,626

Number of rearrangements with t-statistics greater than or equal to 3.888:
105

1-sided p-value from a permutation test of the t-statistic:
105/10626 = .00988

Which test is appropriate?
Independent samples *t* test: exceptionally strong evidence against null ((p=0.00038); Wilcoxon's rank sum test (0.000963, very strong evidence); Permutation test (strong evidence 0.0098); Unequal variance t test (some evidence, 0.038)

| **Slide 51  nCr=24 Choose 4=24!/((24-4)!*4!)** |
| --- |
| |
| NOTES: |
| |
| |
| |
| |
| |

**Matlab solution, normal approximation with ties correction & exact test**

Randomization test: p-value = 0.00988

>>[pvalue,W,U]=Wilcoxranksum(X,Y,0)
pvalue = **9.6335e-004** [With ties correction; identical to SPSS approximation]
W = 84 U = 74

Gallagher's exact Wilcoxon rank sum test in Matlab [algorithms from Hollander & Wolfe (1999)]
>> [pvalue,W,U]=Wilcoxranksum(X,Y,1)
2-tailed pvalue = **0.0038** or **9.4 times Wilcoxon rank sum with ties correction; but 75% of SPSS value, not ties corrected(0.00508)**

| Launch Temperature | Number of O-Ring Incidents |
|---|---|
| Below 65° F | 1 1 1 3 |
| Above 65° F | 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 1 2 |

*EEOS611*

---

**Slide 52  Matlab solution, normal approximation with ties correction & exact test**

NOTES:

---

Test Statistics[b]

| | Incident |
|---|---|
| Mann-Whitney U | 6.000 |
| Wilcoxon W | 216.000 |
| Z | -3.301 |
| Asymp. Sig. (2-tailed) | .00096 |
| Exact Sig. [2*(1-tailed Sig.)] | .00508[a] |

a. Not corrected for ties.

b. Grouping Variable: Temperature

All versions of SPSS, including Version 14:
The exact tests in SPSS are not corrected for ties!

---

**Slide 53  SPSS solution with Wilcoxon Rank sum test, not conservative**

NOTES:

---

**Non-parametric tests are distribution-free, not assumption free**

● No specific distributional assumptions, like normally distributed errors, but all nonparametric tests have some assumptions
● Mann-Whitney U, Underwood 1997, p. 131, "MW/Wilcoxon has nearly identical assumptions to Student's t test"
● **Zar (1999, p. 49) the test is not particularly sensitive to differences in dispersion**
  ‣ Gallagher: Not true in my experience
  ‣ Matlab simulation program available       *EEOS611*

---

**Slide 54  Non-parametric tests are distribution-free, not assumption free**

NOTES:

---

## Randomization doesn't solve problems with unequal variance

- Randomization is often superior to the *t*-distribution for 2-sample problems. It does not remedy the common problems with the *t* distributions though.
- The most common problem with Student's *t* test is the so-called Fisher-Behrens problem: testing for differences in the average if the distributions have different variances
  ‣ This is an open question
  ‣ **Neither** Wilcoxon Rank sum tests **nor** randomization provide a clear solution

*EEOS611*

---

**Slide 55  Randomization doesn't solve problems with unequal variance**

NOTES:

---

## Neither randomization nor permutation tests solve the unequal variance problem

**Manly (1997, p. 141)**
- "The randomization test for the difference in two means can be upset if the samples come from sources that have the same mean, but different variances. This is apparent because the null hypothesis for the randomization test is that the samples come from exactly the same source, which is not true if the variances are not constant"
  ‣ A variety of modifications have been proposed, but all require further study.
- O-ring data may not be a test between mean failure rates!

*EEOS611*

---

**Slide 56  Neither randomization nor permutation tests solve the unequal variance problem**

NOTES:

---

## Gallagher's Matlab Case0401b.m

**Exact tests based on Student's t test**
**Why not just use a 2-sample binomial test?**
>> X = 1 1 1 3; Y = 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 1 1 2

The observed mean difference in O-ring failures was 1.3 incidents per launch with 95% CI of [0.6 2.0]; The t statistic was 3.9 with 22 df and 2-tailed p of **0.00079**;

**Exact tests:** The total number of ways of selecting 4 items from 24 items is 10626; The 1-tailed probability of observing a t statistic greater than observed (3.9) is **0.009881**; The exact 1-tailed p value is: 5/506; The 1-tailed probability of observing a difference greater than observed (1.3) is **0.009881;** The 2-tailed probability of observing a t statistic with abs value greater than observed (|3.9|) is **0.009881...**The exact two-tailed probability for the Wilcoxon rank sum test is **0.003764 The normal approximation for the 2-tailed probability for the Wilcoxon rank sum test is 0.000963**;

The probability of an 0 ring incident if cold (<65F) was 4/4=1.00; The probability of an 0 ring incident in warm (>=65F) was 3/20=0.15; **The two sample binomial test for equal proportions (0.29) of failure has a 2-sided p value of 0.000640**

---

**Slide 57  Gallagher's Matlab Case0401b.m**

NOTES:

## Matlab Case0401b.m

**Summary of conclusions**

The exact two-tailed probability for the Wilcoxon rank sum test is 0.003764 (Can't use: not robust to unequal variances)

The normal approximation for the 2-tailed probability for the Wilcoxon rank sum test is 0.000963 (but this result should NOT be used — it is a large sample approximation)

Binomial test:

The probability of an 0 ring incident if cold (<65F) was 4/4=1.00

The probability of an 0 ring incident in warm (>=65F) was 3/20=0.15

The two-sample binomial test for equal proportions (p=0.29) of failure has a 2-sided p value of 0.000640

---

**Slide 58  Matlab Case0401b.m**

NOTES:

---

## Fligner-Policello test

**Wilcoxon-rank sum test for unequal variances**
**Matlab m.file available, p<3x10$^{-14}$ for O-ring data!**

% Trujillo-Ortiz, A., F. A. Trujillo-Rodriguez, R. Hernandez-Walls, M. A. Fligner and S. Perez-Osuna (2003). FPtest: Non-parametric Fligner-Policello test of two combined random variables with continuous cumulative distribution. A MATLAB file.

%   [WWW document]. URL http://www.mathworks.com/matlabcentral/fileexchange/

%   loadFile.do?objectId=4226&objectType=FILE

% References:

Fligner, M. A. and Policello, G. E. (1981), Robust rank procedure for  the Behrens-Fisher Problem. Journal of the American Statistical Association, 76(373): 162-168.
Hollander, M. and Wolfe, D. (1999), Nonparametric Statistical Methods (2nd ed.).  New York: John Wiley & Sons, Inc. p. 135-139.

---

**Slide 59  Fligner-Policello test**

NOTES:

---

## Asymptotic Power efficiency

**Ratio of sample sizes required to obtain the same p values**

- For normally distributed data, the power efficiency of the Wilcoxon rank sum test is 95.5% of the Student's $t$ test.
  - ‣ For other distributions (e.g., exponential distributions), the power efficiency can be >> 100% (300% for exponential)
    - ▪ Hollander & Wolfe p. 140
- Strengths of Wilcoxon's Rank-sum test
  - ‣ Resistant to outliers
  - ‣ Can handle censored data
- Weakness: generality, determining effect sizes & confidence limits

*EEOS611*

---

**Slide 60  Asymptotic Power efficiency**

NOTES:

| | |
|---|---|
| **Measuring effects sizes**<br>Difficult with nonparametric procedures<br><br>Display 4.8<br><br>Using a rank-sum test to construct a confidence interval for an additive treatment effect; cognitive load study<br><br>*EEOS611* | **Slide 61  Measuring effects sizes**<br><br>NOTES: |
| **Hodges-Lehman estimator for 95% CI**<br>Add a fixed amount to one of the groups: Sleuth's<br><br>Display 4.8<br><br>Using a rank-sum test to construct a confidence interval for an additive treatment effect; cognitive load study<br><br>*EEOS611* | **Slide 62  Hodges-Lehman estimator for 95% CI**<br><br>NOTES: |
| **Unequal variance t test**<br>Welch's t test with Satterthwaite approximation for d.f. | **Slide 63  Unequal variance t test**<br><br>NOTES: |

| | |
|---|---|
|  | **Slide 64  Recall the equal variance t test** |
| | |
| | NOTES: |
| | |
| | |
| | |
| | |
| | |
|  | **Slide 65  Recall Display 2.8, page 40** |
| | |
| | NOTES: |
| | |
| | |
| | |
| | |
| | |
|  | **Slide 66** |
| | |
| | NOTES: |
| | |
| | |
| | |
| | |

| | |
|---|---|
|  | **Slide 67  Welch's t test, p. 97 in Sleuth** |
| | |
| | NOTES: |
| | |
| | |
| | |
| | |
| | |
|  | **Slide 68** |
| | |
| | NOTES: |
| | |
| | |
| | |
| | |
| | |
|  | **Slide 69  Problems with unequal variance t tests (Welch t test)** |
| | |
| | NOTES: |
| | |
| | |
| | |
| | |
| | |

## Example: Stream temperatures

**Unequal variance t test not necessarily conservative**

- Temperatures taken from different portions of a stream:
  - Portion 1: 15.8, 16.9, 17, 17.1, 18, 18.7
    - mean = 17.25, variance = 0.995
  - Portion 2:  18.3, 18.5
    - mean = 18.4,  variance = 0.02
- Obviously the variances are unequal and an equal variance 2-sample t test may be inappropriate
  - Welch [unequal variance]:   t = 2.74 w/ 5 df p = 0.037.
  - Pooled [equal variance]:   t = 1.54 w/ 6 df p = 0.17.
- Why is the equal-variance t-test  giving a lower t-value and a higher p value?

**Slide 70  Example: Stream temperatures**

NOTES:

## The avg temp different: Exact test

**P=3/28, 1-tailed**

There are only 28 ways the 8 temperatures can be arranged into groups of 6 and 2 [8 Choose 2], and in only 3 of these arrangements would the difference in means be equal or greater than the 1.15 ºC difference observed.  These 3 arrangements include the observed data and two others: {18.3, 18.5}, {18.3, 18.7}, {18.5, 18.7}.

P=3/28~0.107

This is the appropriate p value, unless you argue that the variances are different between the 2 portions of stream, but there are too few data to provide strong evidence for this

**Slide 71  The avg temp different: Exact test**

NOTES:

## Levene's test, Section 4.5.3

**Three different tests in the literature**

- Sleuth's Levene's test on page 102-103 is not the same as the Levene's test used by SPSS in Student's *t* test.
  - Sleuth: squared deviations from the mean used as the variables in a *t* test
  - SPSS: absolute values used
    - SPSS:  |observations-mean | used in an F test
    - Other Leven tests |observations-median| used in *t* or F test
- The results can often be quite different
- Levene's tests have largely replaced the $F_{max}$ and Bartlett's tests for equal variance

**Slide 72  Levene's test, Section 4.5.3**

NOTES:

| | |
|---|---|
|  | **Slide 73** |
| | |
| | NOTES: |

| | |
|---|---|
| **4.4 Alternatives to the paired t test**<br><br>Wilcoxon sign-rank and Fisher Sign tests | **Slide 74  4.4 Alternatives to the paired t test** |
| | |
| | NOTES: |

| | |
|---|---|
| **Schizophrenia data**<br>Sleuth p. 99<br><br>Paired *t* test, p=0.006 | **Slide 75  Schizophrenia data** |
| | |
| | NOTES: |

## Slide 76  Wilcoxon signed rank test

NOTES:



## Slide 77  Dealing with tied pairs

NOTES:



## Slide 78  SPSS algorithms, signed rank test

NOTES:

## Fisher's sign test

**Straightforward application of the 2-sample binomial test**

- Given that the probability of a + sign = probability of a minus sign = 0.5,
- What is the probability of observing exactly k positive signs in *n* Bernoulli (binomial) trials
- $P(X=k)= n$ Choose $k * p^k (1-p)^{n-k}$
  - ‣ X has a binomial distribution
  - ‣ Must sum probability for observed value of k, and all more extreme values of k.
- Statistical sleuth provides only the normal approximation to the binomial, but SPSS will provide the exact test for n<30.

**Frequencies**

|  |  | N |
|---|---|---|
| AFFECTED - UNAFFECT | Negative Differences | 14 |
|  | Positive Differences | 1 |
|  | Ties | 0 |
|  | Total | 15 |

a. AFFECTED < UNAFFECT
b. AFFECTED > UNAFFECT
c. AFFECTED = UNAFFECT

**Test Statistics[b]**

|  | AFFECTED - UNAFFECT |
|---|---|
| Exact Sig. (2-tailed) | .0010[a] |

a. Binomial distribution used.
b. Sign Test

*EEOS611*

---

## Sign test in SPSS



*EEOS611*

---

## Conclusions (1 of 2)

**Chapter 4 Alternatives to the t tools**

- Consider using alternatives to the t tools if
  - ‣ The assumptions are grossly violated or
  - ‣ The sample sizes are too small to test distributional assumptions
- Wilcoxon rank sum test
  - ‣ Appropriate for small sample sizes
  - ‣ Appropriate in the presence of outliers
  - ‣ Ties are not a problem if the ties-correction used
  - ‣ Not appropriate for samples with unequal variances (try Fligner-Policello if the sample sizes are large)

*EEOS611*

---

### Slide 79  Fisher's sign test

NOTES:

### Slide 80  Sign test in SPSS

NOTES:

### Slide 81  Conclusions (1 of 2)

NOTES:

## Conclusions (2 of 2)

### Chapter 4 Alternatives to the t tools

- Permutation test
  - Appropriate for small sample sizes, when the Student's *t* distribution might not be appropriate
  - Does not protect against the problem of unequal variances (the Fisher-Behrens problem)
- Paired data: tests based on ranks
  - Wilcoxon signed rank test: high power efficiency
  - Sign test, simple application of the 1-sample binomial

*EEOS611*

**Slide 82  Conclusions (2 of 2)**

NOTES: