Statistics We Will Consider				
DV 🕨	Categorical	Parametric Interval/ND	Nonparametric Ordinal/~ND	
univariate stats	mode, #cats	mean, std	median, IQR	
univariate tests	gof X <sup>2</sup>	1-grp t-test	1-grp Mdn test	
association	X <sup>2</sup>	Pearson's r	Spearman's r Kendall's Tau	
2 bg	X <sup>2</sup>	t- / F-test	M-W K-W Mdn	
k bg	X <sup>2</sup>	F-test	K-W Mdn	
2wg	McNem Crn's	t- / F-test	Wil's Fried's	
kwg	Crn's	F-test	Friedman's	
M-W Mann-Wh K-W Kruskal-W Mdn Median Te	nitney U-Test Wil's Vallis Test st McN	Wilcoxin's Test Fr em McNemar's X <sup>2</sup>	ied's Friedman's F-test Crn's – Cochran's Test	
	DV → univariate stats univariate tests association 2 bg k bg 2wg kwg M-W Mann-Wi K-W Kruskal-W Mdn Median Te	DV       ◆       Categorical         univariate stats       mode, #cats         univariate tests       gof X²         association       X²         2 bg       X²         k bg       X²         2wg       McNem Crn's         M-W Mann-Whitney U-Test       Wil's         K-W Kruskal-Wallis Test       McN	Statistics We Will Consid         DV       Categorical       Parametric         univariate stats       mode, #cats       mean, std         univariate tests       gof X <sup>2</sup> 1-grp t-test         association       X <sup>2</sup> Pearson's r         2 bg       X <sup>2</sup> F-test         k bg       X <sup>2</sup> F-test         2wg       McNem Crn's       t- / F-test         kwg       Crn's       F-test         M-W Mann-Whitney U-Test       Wil's Wilcoxin's Test       Fr         Mdn Median Test       McNem McNemar's X <sup>2</sup>	

### Statistical Tests for BG Designs w/ qualitative variables

Pearson's X<sup>2</sup>

$$X^2 = \sum \frac{(of - ef)^2}{ef}$$

Can be 2x2 or kxk – depending upon the number of categories of the qualitative outcome variable

- H0: Populations represented by the design conditions have the same distribution across conditions/categories of the outcome variable
- degrees of freedom df = (#colums 1) \* (#rows 1)
- Range of values 0 to  $\infty$
- Reject Ho: If  $X^2_{obtained} > X^2_{critical}$



The expected frequency for each cell is computed assuming that the H0: is true – that there is no relationship between the row and column variables.

If so, the frequency of each cell can be computed from the frequency of the associated rows & columns.



(76\*86)/154

(78\*86)/154

86

76

78

154



df = (2-1) \* (2-1) = 1



 $X^{2}_{1.05} = 3.84$ 

 $X_{1,.01}^2 = 6.63$ 

p = .0002 using online p-value calculator

So, we would reject H0: and conclude that the two groups have different distributions of responses of the qualitative DV.

Parametric tests for BG Designs using ND/Int variables

#### t-tests

• H0: Populations represented by the IV conditions have the same mean DV.

Row 1

Row 2

(76\*68)/154

(78\*68)/154

68

- degrees of freedom df = N 2
- Range of values  $-\infty$  to  $\infty$
- Reject Ho: If  $|t_{obtained}| > t_{critical}$
- Assumptions
  - data are measured on an interval scale
  - DV values from both groups come from ND with equal STD

# ANOVA

- H0: Populations represented by the IV conditions have the same mean DV.
- degrees of freedom df numerator = k-1, denominator = N k
- Range of values 0 to  $\infty$
- Reject Ho: If F<sub>obtained</sub> > F<sub>critical</sub>
- Assumptions
  - data are measured on an interval scale
  - DV values from both groups come from ND with equal STD

The nonparametric BG models we will examine, and the parametric BG models with which they are most similar...

# 2-BG Comparisons

Mann-Whitney U test

between groups t-test

## 2- or k-BG Comparisons

Kruskal-Wallis test

Median test

between groups ANOVA

between groups ANOVA

As with parametric tests, the k-group nonparametric tests can be used with 2 or k-groups.

Let's start with a review of applying a between groups t-test

Here are the data from such a design :

Qual variable is whether or not subject has a 2-5 year old Quant variable is "liking rating of Barney" (1-10 scale)

No Too	dlertoddle	r 1+ Tod	dlers	_
s1	2	s3	6	
s2	4	s5	8	
s4	6	s6	9	
s8	7	s7	10	
M =	4.75	M =	8.25	

The BG t-test would be used to compare these group means.

When we perform this t-test ...

As you know, the H0: is that the two groups have the same mean on the quantitative DV, but we also ...

- Assume that the quantitative variable is measured on a interval scale -- that the difference between the ratings of "2" and "4" mean the same thing as the difference between the ratings of "8" and "6".
- 2. Assume that the quant variable is normally distributed.

3. Assume that the two samples have the same variability (homogeneity of variance assumption)

Given these assumptions, we can use a t-test tp assess the

H0:  $M_1 = M_2$ 

If we want to "avoid" these first two assumptions, we can apply the No Toddlestoddler 1+ Toddlers Mann-Whitney U-test All the values are rating ranks rating ranks ranked at once --2 3.5 s1 1 s3 6 ignoring which The test does not depend upon the interval properties of the data. only their ordinal properties -- and so we will convert the values to condition each "S" s2 8 4 2 s5 6 ranks was in. lower scores have lower ranks, and vice versa 9 7 Notice the group s4 6 3.5 s6 • e.a. #1 values 10 11 13 14 16 with the higher ranks 1 2 3 4 5 5 values has the s8 7 s7 10 8 • Tied values given the "average rank" of all scores with that value higher summed • e.g. #2 values 10 12 12 13 16  $\Sigma = 11.5$  $\Sigma = 24.5$ ranks ranks 12.5 2.5 4 5 • e.g., #3 values 9 12 13 13 13 The "U" statistic is computed from the summed ranks. U=0 when ranks 1 2 4 4 4 the summed ranks for the two groups are the same (H0:) There are two different "versions" of the H0: for the Mann-Whitney U-test, depending upon which text you read. The "older" version reads: H0: The samples represent populations with the same distributions of scores. Under this H0:, we might find a significant U because the samples from the two populations differ in terms of their: centers (medians - with rank data) variability or spread shape or skewness This is a very "general" H0: and rejecting it provides little info. Also, this H0: is not strongly parallel to that of the t-test (which is specifically about mean differences)

Nonparametric tests for BG Designs using ~ND/~Int variables

Preparing these data for analysis as ranks...

Over time, "another" H0: has emerged, and is more commonly seen in textbooks today:

H0: The two samples represent populations with the same median (assuming these populations have distributions with identical variability and shape).

You can see that this H0:

- increases the specificity of the H0: by making assumptions (That's how it works - another one of those "trade-offs")
- is more parallel to the H0: of the t-test (both are about "centers")
- has essentially the same distribution assumptions as the t-test (equal variability and shape)

Finally, there are two "forms" of the Mann-Whitney U-test:

With smaller samples (n < 20 for both groups)

- $\mbox{-}$  compare the summed ranks to the two groups to compute the test statistic -- U
- -Compare the  $W_{\text{obtained}}$  with a  $W_{\text{critical}}$  that is determined based on the sample size

With larger samples (n > 20)

- with these larger samples the distribution of U-obtained values approximates a normal distribution
- a Z-test is used to compare the Uobtained with the Ucritical
- the  $Z_{obtained}$  is compared to a critical value of 1.96 (p = .05)

Nonparametric tests for BG Designs using ~ND/~Int variables

The Kruskal- Wallis test

- applies this same basic idea as the Mann-Whitney Utest (comparing summed ranks)
- can be used to compare any number of groups.
- DV values are converted to rankings
  - ignoring group membership
  - assigning average rank values to tied scores
- Score ranks are summed within each group and used to compute a summary statistic "H", which is compared to a critical value obtained from a X<sup>2</sup> distribution to test H0:
  - groups with higher values will have higher summed ranks
  - if the groups have about the same values, they will have about the same summed ranks

H0: has same two "versions" as Mann-Whitney U-test	Nonparametric tests for BG Designs using ~ND/~Int variables
<ul> <li>groups represent populations with same score distributions</li> </ul>	Median Test also for comparing 2 or multiple groups
<ul> <li>groups represent pops with same median (assuming these populations have distributions with identical variability and shape).</li> </ul>	The intent of this test was to compare the medians of the groups, without the "distributions are equivalent" assumptions of the Mann-Whitney and Kruskal-Wallis tests
<ul> <li>Rejecting H0: tells only that there is some pattern of distribution/median difference among the groups</li> <li>specifying this pattern requires pairwise K-W follow-up analyses</li> <li>Bonferroni correction p<sub>critical</sub> = (.05 / # pairwise comps)</li> </ul>	<ul> <li>This was done in a very creative way</li> <li>compute the grand median (ignoring group membership)</li> <li>for each group, determine which members have scores above the grand median, and which have scores below the grand median</li> </ul>
Assemble the information into a contingency table • Perform a Pearson's (contingency table) X <sup>2</sup> to test for a pattern of median differences (pairwise follow-ups) • Please note: The median test has substantially less power than the Kruskal-Wallis test for the same sample size e.g., Mdn <sub>1</sub> = Mdn <sub>2</sub> = Mdn <sub>3</sub> G <sub>1</sub> G <sub>2</sub> G <sub>3</sub> • $G_1$ G <sub>2</sub> G <sub>3</sub> • $C_1$ G <sub>2</sub> G <sub>3</sub> • $C_2$ G <sub>3</sub> • $C_1$ G <sub>2</sub> G <sub>3</sub> • $C_2$ G <sub>3</sub> • $C_1$ G <sub>2</sub> G <sub>3</sub> • $C_1$ G <sub>2</sub> G <sub>3</sub> • $C_2$ G <sub>3</sub> • $C_1$ G <sub>2</sub> G <sub>3</sub> • $C_2$ G <sub>3</sub> • $C_1$ G <sub>2</sub> G <sub>3</sub> • $C_1$ G <sub>2</sub> G <sub>3</sub> • $C_2$ C <sub>3</sub>	

Repeated measures designs...

There are two major kinds of these designs:

- 1) same cases measured on the same variable at different times or under different conditions
  - pre-test vs. post-test scores of clients receiving therapy
  - performance scores under feedback vs. no feedback conds
  - % who "pass" before versus after remedial training
- 2) same cases measured at one time under one condition, using different (yet comparable) measures
  - comparing math and reading scores (both T-scores, with mean=50 and std=10)
  - number of "omissions" (words left out) and "intrusions" (words that shouldn't have been included) in a word recall task
  - % who "pass" using two different tests

Repeated measures designs...

There is really a third related kind of design:

- 3) non-independent groups of cases measured on the same variable at different times or under different conditions
  - matched-groups designs
  - snow-ball sampling over time

Statistically speaking, groups-comparisons analyses divide into 2 kinds"

- independent groups designs  $\rightarrow$  Between Groups designs
- dependent groups designs  $\rightarrow$  within-groups & Matchedgroups designs

For all dependent groups designs, the non-independence of the groups allows the separation of variance due to "differences among people" from variance due to "unknown causes" (error or residual variance)

For repeated measures designs (especially of the first 2 kinds), there are two different types of research hypotheses or questions that might be posed...

1) Do the measures have different means (dif resp dist for qual DVs)

- are post-test scores higher than pre-test scores?
- is performance better with feedback than without it?
- are reading scores higher than math scores?
- are there more omissions than intrusions?

2) Are the measures associated?

- are the folx with the highest pre-test scores also the ones with the highest post-test scores?
- is performance with feedback predictable based on performance without feedback?
- are math scores and reading scores correlated?
- do participants who make more omissions also tend to make more intrusions?

So, taken together there are four "kinds of" repeated measures analyses. Each is jointly determined by the type of design and the type of research hypothesis/question. Like this...

					Consider the difference between the following examples of repeated measures designs using a qualitative (binary) response or outcome variable
	Т	ype of Hypoth	esis/Quest	ion	<ul> <li>The same % of students will be identified as needing remedial instruction at the beginning and end of the semester (dif times)</li> </ul>
Type of Design	mean	difference	assoc	ciation	•The same students will be identified as needing remedial instruction at the end of the semester as at the beginning (dif times)
Different times or situations	pre-tes	t < post-test	pre-test	& post-test	• The same % of students will be identified as needing remedial instruction based on teacher evaluations as based on a standardized test (dif measures)
Different measure	s math	< spelling	math 8	spelling	<ul> <li>The same students will be identified as needing remedial instruction based on teacher evaluations as based on a standardized test (dif measures)</li> </ul>
					So, we have to expand our thinking to include 8 situations
So, for repeated r "situations" and th	neasures d e statistic t	lesigns, here a to use for eacl	are the ana	lytic	
	Т	ype of Question	on/Hypothe	esis	
	Q	uant Vars	Qua	al Vars	
Type of Design	mean dif	assoc	% dif^	pattern^*	
Different times or situations	wg t/F-test	Pearson's r	Cochrans	McNemar's X <sup>2</sup>	
Different measures	wg t/F-test	Pearson's r	Cochran's	McNemar's X <sup>2</sup>	

But... All the examples so far have used quantitative variables.

Qualitative variables could be used with each type of repeated measures design (dif times vs. dif measures)

 $^{\mbox{\sc h}}$  Cochran's and McNemar's are for use only with binary variables

\* McNemar's looks at patterns of classification disagreements

### Statistical Tests for WG Designs w/ qualitative variables

## McNemar's test

Of all these tests, McNemar's has the most specific application...

- are two qualitative variable related -- Pearson's X<sup>2</sup>
- do groups have differences on a qual variable -- Pearson's X<sup>2</sup>
- · does a group change % on a binary variable -- Cochran's

• is the the relationship between the variables revealed by an asymmetrical pattern of "disagreements" –McNemar's



Cochran's Q-test - can be applied to 2 or k-groups

The simplest "qualitative variable" situation is when the variable is binary. Then "changes in response distribution" becomes the much simple "changes in %".

Begin the computation of Q by arranging the data with each case on a separate row. 1 = pass 0 = fail

S1	pretest 0	posttest 1	retention 1	L 2	L² 4	
S2	0	1	0	1	1	
S2	0	0	1	1	1	
S2	1	1	1	3	9	
S5	0	0	1	1	1	
G	1	3	4	+		
G <sup>2</sup>	1	9	16			
Com colur	pute the su nn (G) and	m for each it's square (G²)	1	Compute row (L) a	e the s and its	um for each square (L²)



Compare the obtained X<sup>2</sup> with X<sup>2</sup>  $_{1,.05}$  = 3.84. We would reject H0: and conclude that there is a relationship between what performance on the paper test and performance on the computer test & that more uniquely fail the paper test than uniquely fail the computer test.

$\begin{array}{c ccccccccccccccccccccccccccccccccccc$	Parametric tests for WG Designs using ND/Int variables t-tests • H0: Populations represented by the IV conditions have the same mean DV. • degrees of freedom df = N - 1 • Range of values -∞ to ∞ • Reject H0: If   t <sub>obtained</sub>   > t <sub>critical</sub> • Assumptions • data are measured on an interval scale • DV values from both groups come from ND & have equal STDs
$(k-1)^*[(k * \Sigma G^2) - (\Sigma G)^2]$ (3-1)* $[(3^*(1+9+16)) - (1+3+4)^2]$	ANOVA • H0: Populations represented by the IV conditions have the same mean DV
$\mathbf{Q} = \frac{1}{(k + \Sigma L) - \Sigma L^2} = \frac{1}{(3 + (2+1+1+3+1)) - (4+1+1+9+1)} = 3.0$	• degrees of freedom df numerator = $k-1$ , denominator = $N - k$
Q is compared to X <sup>2</sup> critical based on df = k-1 $X_{2,.05}^2 = 7.81$ So we would retain H0: of no % difference across the design conditions.	<ul> <li>Range of values 0 to ∞</li> <li>Reject Ho: If F<sub>obtained</sub> &gt; F<sub>critical</sub></li> <li>Assumptions         <ul> <li>data are measured on an interval scale</li> <li>DV values from both groups come from ND with equal STD</li> <li>for k &gt; 2 – data from any pair of conditions are equally correlated</li> </ul> </li> </ul>
Nonparametric tests for WG Designs using ~ND/~Int variables • within-subjects design - same subjects giving data under each of two or more conditions	
<ul> <li>comparison of two or more "comparable" variables same subjects giving data on two variables (same/dif time)</li> </ul>	
<ul> <li>matched-groups design matched groups of two or more members, each in one of the conditions</li> </ul>	
The nonparametric RM models we will examine and their closes parametric RM counterparts	t
2-WG Comparisons Wilcoxin's Test dependent t-test	
2- or k-WG Comparisons Friedman's ANOVA dependent ANOVA	

Let's start with a review of applying a within groups t-test

Here are the data from such a design : IV is Before vs. After the child "discovers" Barney (and watches it incessantly, exposing you to it as well) so..

1st Quant variable is 1-10 rating "before" discovery 2nd Quant variable is 1-10 rating "after discovery"

Before	After	Difference
s1 2	s1 6	-4
s2 4	s2 8	-4
s3 6	s3 9	-3
s4 7	s4 10	-3
M = 4.75	M = 8.25	M <sub>d</sub> = -3.5

A WG t-test can be computed as a single-sample t-test using the differences between an individual's scores from the 2 design conditions.

• Rejecting the H0:  $M_d=0$ , is rejecting the H0:  $M_{before} = M_{after}$ 

• other formulas exist

When using a WG t-test (no matter what computational form the assumption of interval measurement properties is even "more assuming" than for the BG design. We assume ...

- that each person's ratings are equally spaced -- that the difference between ratings given by S1 of "3" and "5" mean the same thing as the difference between their ratings of "8" and "10" ???
- that different person's rating are equally spaced -- that the difference between ratings given by S1 of "3" and "5" mean the same thing as the difference between ratings of "8" and "10" given by S2 ???

Nonparametric tests for WG Designs using ~ND/~Int variables

# Wilcoxin's Test

If we want to avoid some assumptions, we can apply a nonparametric test. To do that we ...

- Compute the differences between each person's scores
- Determine the "signed ranks" of the differences
- Compute the summary statistic W from the signed ranks

Before After		Difference	Signed Ranks		
s1	2	s1	5	3	2.5
s2	4	s2	8	4	4
s3	6	s3	9	3	2.5
s4	9	s4	7	-2	-1

The "W" statistic is computed from the signed ranks. W=0 when the signed ranks for the two groups are the same (H0:)

There are two different "versions" of the H0: for the Wilcoxin's test, depending upon which text you read.

The "older" version reads:

- H0: The two sets of scores represent a population with the same distribution of scores under the two conditions.
- Under this H0:, we might find a significant U because the samples from the two situations differ in terms of their:
  - centers (medians with rank data)
  - · variability or spread
  - shape or skewness

This is a very "general" H0: and rejecting it provides little info.

Also, this H0: is not strongly parallel to that of the t-test (that is specifically about mean differences)

Over time, "another" H0: has emerged, and is more commonly seen in textbooks today:

H0: The two sets of scores represent a population with the same median under the two conditions (assuming these populations have distributions with identical variability and shape).

You can see that this H0:

- increases the specificity of the H0: by making assumptions (That's how it works - another one of those "trade-offs")
- is more parallel to the H0: of the t-test (both are about "centers")
- has essentially the same distribution assumptions as the t-test (equal variability and shape)

Finally, there are also "forms" of the Wilcoxin's Test:

With smaller samples (N < 10-50 depending upon the source ??)

- Compare the  $W_{\text{obtained}}$  with a  $W_{\text{critical}}$  that is determined based on the sample size

With larger samples (N > 10-50)

- with these larger samples the distribution of Uobtained values approximates a normal distribution
- a Z-test is used to compare the Uobtained with the Ucritical
- the  $Z_{obtained}$  is compared to a critical value of 1.96 (p = .05)

You should notice considerable similarity between the Mann-Whitney U-test and the Wilcoxin -- in fact, there are BG and RM versions of each -- so be sure to ask the "version" whenever you hear about one of these tests. Nonparametric tests for WG Designs using ~ND/~Int variables

Friedman's test applies this same basic idea (comparing ranks), but can be used to compare any number of groups.

- Each subject's DV values are converted to rankings (across IV conditions)
- Score ranks are summed within each IV Condition and used to compute a summary statistic "F", which is compared to a critical value to test H0:
- E.g., -- more of Barney . . . (from different "stages of exposure")

	Before		After 6 months		After 12 months		S
	DV	rank	DV	rank	DV	rank	
S1	3	1	7	3	5	2	
S2	5	1	9	3	6	2	
S3	4	2	6	3	2	1	
S4	3	1	6	2	9	3	

- H0: has same two "versions" as the other nonparametric tests
  - DVs from populations with same score distributions
  - DVs from populations with same median (assuming ...)
- Rejecting H0: requires pairwise follow-up analyses

• Bonferroni correction -- p<sub>critical</sub> = (.05 / # pairwise comps)

- Finally, there are also "forms" of Friedman's Test:
  - With smaller samples (k < 6 & N < 14)
    - Compare the F<sub>obtained</sub> with a F<sub>critical</sub> that is determined based on the sample size & number of conditions
  - With larger samples (k > 6 or N > 14)
    - the  $\mathsf{F}_{obtained}$  is compared to a  $X^2_{critical}$  value