

# Statistical Inference: matched pairs and normal quantile plot

# What to do if normality fails

- If normality fails (for one or both of the groups), what do we do then?
- Again, can compare medians: use the thought process of the sign test, which does not depend on normality and is not damaged by outliers.
- A suitable test called Mood's median test.
- Before we get to that, a diversion.

# The chi-squared test for independence

Suppose we want to know whether people are in favour of having daylight savings time all year round. We ask 20 males and 20 females whether they each agree with having DST all year round ("yes") or not ("no"). Some of the data:

```
my_url="http://www.utsc.utoronto.ca/~butler/c32/dst.txt"
dst=read_delim(my_url," ")
dst %>% sample_n(5) # randomly sample 5 rows
```

gender	agree
male	yes
male	yes
male	no
male	no
female	no

## ... continued

Count up individuals in each category combination, and arrange in contingency table:

```
tab=with(dst, table(gender, agree))  
tab
```

```
##           agree  
## gender    no  yes  
##  female  11   9  
##   male   3  17
```

- Most of the males say “yes”, but the females are about evenly split.
- Looks like males more likely to say “yes”, ie. an association between gender and agreement.
- Test an  $H_0$  of “no association” (“independence”) vs. alternative that there is really some association.
- Done with `chisq.test`.

## ...And finally

```
chisq.test(tab,correct=F)
```

```
##  
##  Pearson's Chi-squared test  
##  
## data:  tab  
## X-squared = 7.033, df = 1, p-value = 0.008002
```

- Reject null hypothesis of no association
- therefore there is a difference in rates of agreement between (all) males and females (or that gender and agreement are associated).
- Without `correct=F` uses “Yates correction”; this way, should give same answers as calculated by hand (if you know how).

# Mood's median test

- Before our diversion, we wanted to compare medians of two groups.
- Recall sign test: count number of values above and below something (there, hypothesized median).
- Idea of Mood's median test:
  - Work out the median of all the data, regardless of group ("grand median").
  - Count how many data values in each group are above/below this grand median.
  - Make contingency table of group vs. above/below.
  - Test for association.
- If group medians equal, each group should have about half its observations above/below grand median. If not, one group will be mostly above grand median and other below.

# Mood's median test for reading data

- Find overall median score:

```
(kids %>% summarize(med=median(score)) %>% pull(med) -> m)
```

```
## [1] 47
```

- Make table of above/below vs. group:

```
tab=with(kids, table(group, score>m))  
tab
```

```
##
```

```
## group FALSE TRUE
```

```
##      c      15      8
```

```
##      t       7     14
```

- Treatment group scores mostly above median, control group scores mostly below, as expected.

# The test

- Do chi-squared test:

```
chisq.test(tab,correct=F)
```

```
##
```

```
## Pearson's Chi-squared test
```

```
##
```

```
## data:  tab
```

```
## X-squared = 4.4638, df = 1, p-value = 0.03462
```

- This test actually two-sided (tests for any association).
- Here want to test that new reading method *better* (one-sided).
- Most of treatment children above overall median, so do 1-sided test by halving P-value to get 0.017.
- This way too, children do better at learning to read using the new method.



# Or by smmr

- `median_test` does the whole thing:

```
median_test(kids,score,group)
```

```
## $table
##      above
## group above below
##      c      8      15
##      t     14       7
##
## $test
##      what      value
## 1 statistic 4.46376812
## 2          df 1.00000000
## 3    P-value 0.03462105
```

- P-value again two-sided.

# Comments

- P-value 0.013 for (1-sided)  $t$ -test, 0.017 for (1-sided) Mood median test.
- Like the sign test, Mood's median test doesn't use the data very efficiently (only, is each value above or below grand median).
- Thus, if we can justify doing  $t$ -test, we should do it. This is the case here.
- The  $t$ -test will usually give smaller P-value because it uses the data more efficiently.
- The time to use Mood's median test is if we are definitely unhappy with the normality assumption (and thus the  $t$ -test P-value is not to be trusted).

# Jumping rats

- Link between exercise and healthy bones (many studies).
- Exercise stresses bones and causes them to get stronger.
- Study (Purdue): effect of jumping on bone density of growing rats.
- 30 rats, randomly assigned to 1 of 3 treatments:
  - No jumping (control)
  - Low-jump treatment (30 cm)
  - High-jump treatment (60 cm)
- 8 weeks, 10 jumps/day, 5 days/week.
- Bone density of rats ( $\text{mg}/\text{cm}^3$ ) measured at end.
- See whether larger amount of exercise (jumping) went with higher bone density.
- Random assignment: rats in each group similar in all important ways.
- So entitled to draw conclusions about cause and effect.

# Reading the data

Values separated by spaces:

```
my_url="http://www.utsc.utoronto.ca/~butler/c32/jumping.txt"
rats=read_delim(my_url, " ")
```

```
## Parsed with column specification:
```

```
## cols(
```

```
##   group = col_character(),
```

```
##   density = col_double()
```

```
## )
```

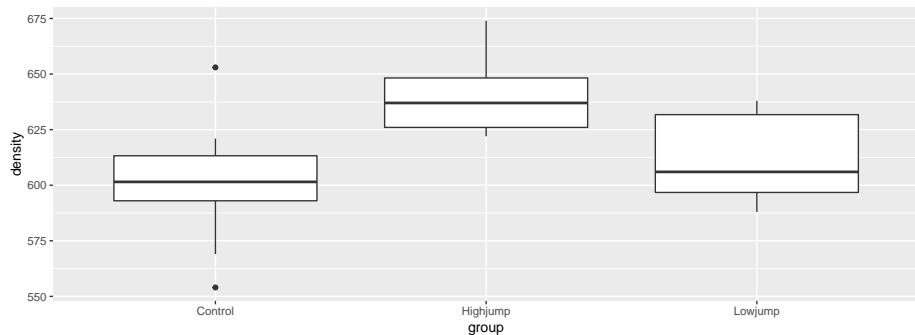
## The data (some random rows)

```
rats %>% sample_n(12)
```

group	density
Highjump	650
Control	554
Control	614
Control	621
Highjump	643
Lowjump	594
Control	600
Lowjump	596
Control	593
Lowjump	638
Lowjump	588
Highjump	626

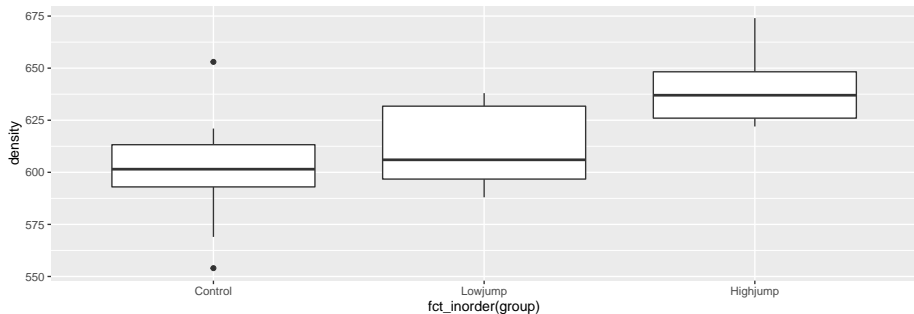
# Boxplots

```
ggplot(rats, aes(y=density, x=group)) + geom_boxplot()
```



## Or, arranging groups in data (logical) order

```
ggplot(rats, aes(y=density, x=fct_inorder(group))) +  
geom_boxplot()
```



# Analysis of Variance

- Comparing  $> 2$  groups of independent observations (each rat only does one amount of jumping).
- Standard procedure: analysis of variance (ANOVA).
- Null hypothesis: all groups have same mean.
- Alternative: “not all means the same”, at least one is different from others.



# Testing: ANOVA in R

```
rats.aov=aov(density~group,data=rats)
summary(rats.aov)
```

```
##              Df Sum Sq Mean Sq F value Pr(>F)
## group         2    7434     3717   7.978 0.0019 **
## Residuals    27   12579       466
## ---
## Signif. codes:
## 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

- Usual ANOVA table, small P-value: significant result.
- Conclude that the mean bone densities are not all equal.
- Reject null, but not very useful finding.

# Which groups are different from which?

- ANOVA really only answers half our questions: it says “there are differences”, but doesn’t tell us which groups different.
- One possibility (not the best): compare all possible pairs of groups, via two-sample t.
- First pick out each group:

```
rats %>% filter(group=="Control") -> controls  
rats %>% filter(group=="Lowjump") -> lows  
rats %>% filter(group=="Highjump") -> highs
```

## Control vs. low

```
t.test(controls$density, lows$density)
```

```
##  
## Welch Two Sample t-test  
##  
## data: controls$density and lows$density  
## t = -1.0761, df = 16.191, p-value = 0.2977  
## alternative hypothesis: true difference in means is not equal to 0  
## 95 percent confidence interval:  
## -33.83725 11.03725  
## sample estimates:  
## mean of x mean of y  
## 601.1 612.5
```

No sig. difference here.

# Control vs. high

```
t.test(controls$density, highs$density)
```

```
##  
## Welch Two Sample t-test  
##  
## data: controls$density and highs$density  
## t = -3.7155, df = 14.831, p-value = 0.002109  
## alternative hypothesis: true difference in means is not equal to 0  
## 95 percent confidence interval:  
## -59.19139 -16.00861  
## sample estimates:  
## mean of x mean of y  
## 601.1 638.7
```

These are different.

## Low vs. high

```
t.test( lows$density, highs$density)
```

```
##  
##  Welch Two Sample t-test  
##  
## data:  lows$density and highs$density  
## t = -3.2523, df = 17.597, p-value = 0.004525  
## alternative hypothesis: true difference in means is not equal to 0  
## 95 percent confidence interval:  
##  -43.15242  -9.24758  
## sample estimates:  
## mean of x mean of y  
##      612.5      638.7
```

These are different too.

# But...

- We just did 3 tests instead of 1.
- So we have given ourselves 3 chances to reject  $H_0$  : all means equal, instead of 1.
- Thus  $\alpha$  for this combined test is not 0.05.

# John W. Tukey



- American statistician, 1915–2000
- Big fan of exploratory data analysis
- Invented boxplot
- Invented "honestly significant differences"
- Invented jackknife estimation
- Coined computing term "bit"
- Co-inventor of Fast Fourier Transform

# Honestly Significant Differences

- Compare several groups with one test, telling you which groups differ from which.
- Idea: if all population means equal, find distribution of highest sample mean minus lowest sample mean.
- Any means unusually different compared to that declared significantly different.



# Tukey on rat data

```
rats.aov=aov(density~group,data=rats)
TukeyHSD(rats.aov)
```

```
##    Tukey multiple comparisons of means
##      95% family-wise confidence level
##
## Fit: aov(formula = density ~ group, data = rats)
##
## $group
```

		diff	lwr	upr	p adj
## Highjump-Control	37.6	13.66604	61.533957	0.0016388	
## Lowjump-Control	11.4	-12.53396	35.333957	0.4744032	
## Lowjump-Highjump	-26.2	-50.13396	-2.266043	0.0297843	

- Again conclude that bone density for highjump group significantly higher than for other two groups.

# Why Tukey's procedure better than all t-tests

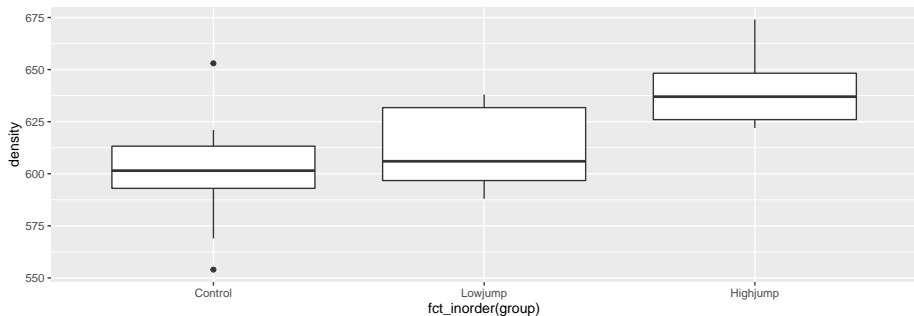
Look at P-values for the two tests:

Comparison	Tukey	t-tests
-----		
Highjump-Control	0.0016	0.0021
Lowjump-Control	0.4744	0.2977
Lowjump-Highjump	0.0298	0.0045

- Tukey P-values (mostly) higher.
- Proper adjustment for doing three t-tests at once, not just one in isolation.
- lowjump-highjump comparison would no longer be significant at  $\alpha = 0.01$ .

# Checking assumptions

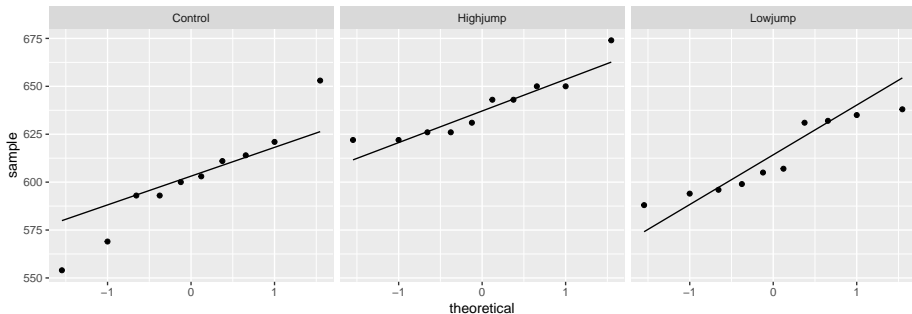
```
ggplot(rats, aes(y=density, x=fct_inorder(group))) +  
geom_boxplot()
```



Assumptions: - Normally distributed data within each group - with equal group SDs.

# Normal quantile plots by group

```
ggplot(rats, aes(sample = density)) + stat_qq() +  
  stat_qq_line() + facet_wrap(~ group)
```



# The assumptions

- Normally-distributed data within each group
- Equal group SDs. These are shaky here because:
  - control group has outliers
  - highjump group appears to have less spread than others. Possible remedies (in general):
- Transformation of response (usually works best when SD increases with mean)
- If normality OK but equal spreads not, can use Welch ANOVA. (Regular ANOVA like pooled t-test; Welch ANOVA like Welch-Satterthwaite t-test.)
- Can also use Mood's Median Test (see over). This works for any number of groups.

# Mood's median test 1/4

- Find median of all bone densities, regardless of group:

```
(rats %>% summarize(med = median(density)) %>% pull(med) -> m)
```

```
## [1] 621.5
```

- Count up how many observations in each group above or below overall median:

```
tab = with(rats, table(group, density > m))  
tab
```

```
##  
## group      FALSE TRUE  
## Control      9     1  
## Highjump     0     10  
## Lowjump      6     4
```

## Mood's median test 2/4

```
tab
```

```
##
```

```
## group      FALSE TRUE
```

```
##   Control      9    1
```

```
##   Highjump     0   10
```

```
##   Lowjump      6    4
```

- All Highjump obs above overall median.
- Most Control obs below overall median.
- Suggests medians differ by group.

## Mood's median test 3/4

- Test whether association between group and being above/below overall median significant using chi-squared test for association:

```
chisq.test(tab,correct=F)
```

```
##
```

```
## Pearson's Chi-squared test
```

```
##
```

```
## data:  tab
```

```
## X-squared = 16.8, df = 2, p-value = 0.0002249
```

- Very small P-value says that being above/below overall median depends on group.
- That is, groups do not all have same median.



## Mood's median test 4/4

Or with `median_test` from `smmr`, same as before.

```
median_test(rats,density,group)
```

```
## $table
##           above
## group      above below
##   Control      1     9
##   Highjump     10     0
##   Lowjump      4     6
##
## $test
##           what           value
## 1 statistic 1.680000e+01
## 2          df 2.000000e+00
## 3    P-value 2.248673e-04
```

# Comments

- No doubt that medians differ between groups (not all same).
- This test is equivalent of  $F$ -test, not of Tukey.
- To determine which groups differ from which, can compare all possible pairs of groups via (2-sample) Mood's median tests, then adjust P-values by multiplying by number of 2-sample Mood tests done (Bonferroni):

```
pairwise_median_test(rats,density,group)
```

g1	g2	p_value	adj_p_value
Control	Highjump	0.0001478	0.0004434
Control	Lowjump	0.3710934	1.0000000
Highjump	Lowjump	0.3710934	1.0000000

- Now, lowjump-highjump difference no longer significant.

# Welch ANOVA

- For these data, Mood's median test probably best because we doubt both normality and equal spreads.
- When normality OK but spreads differ, Welch ANOVA way to go.
- Welch ANOVA done by `oneway.test` as shown (for illustration):

```
oneway.test(density~group,data=rats)
```

```
##  
## One-way analysis of means (not assuming  
## equal variances)  
##  
## data: density and group  
## F = 8.8164, num df = 2.000, denom df =  
## 17.405, p-value = 0.002268
```

- P-value very similar, as expected.
- Appropriate Tukey-equivalent here called Games-Howell.

# Games-Howell

- Lives in package PMCMRplus (also userfriendlyscience). Install first.

```
library(PMCMRplus)
```

```
gamesHowellTest(density~factor(group),data=rats)
```

```
##  
## Pairwise comparisons using Games-Howell test  
## data: density by factor(group)  
##  
##           Control Highjump  
## Highjump 0.0056  -  
## Lowjump  0.5417  0.0120  
##  
## P value adjustment method: none
```