

ORIGIN ≡ 1

Hotelling's T² Test for a Paired Design

Just as with the univariate t-test, Hotelling's T² test for paired data is fundamentally identical to the corresponding one-population version of the test. All one does is use the difference D matrix between matched observations X and Y and treat this as a single population. As with the univariate case, paired tests have greater power than non-paired two-sample T² tests because variance between individuals, unrelated to treatments, is better controlled. The example comes from RA Johnson & DW Wichern (JW) *Applied Multivariate Statistical Analysis 4th Edition 1998*. It is useful in demonstrating that rejection of the null hypothesis is possible even when all confidence intervals include the vector of zero difference. In these cases, the multivariate T² confidence ellipsoid, identical with the T² test, should be preferred.

Effluent Data JW Table 6-1 p. 275)

Samples taken and split for comparison of technique between:

Read Data:

X = variables BOD & SS measured in a Commercial Lab.
Y = variables BOD & SS measured in a State Lab.

```
M := READPRN("c:\DATA\Multivariate\T6-1.DAT")
```

```
n := rows(M)      n = 11
```

```
X := submatrix(M, 1, n, 1, 2)
```

```
Y := submatrix(M, 1, n, 3, 4)
```

```
p := cols(X)      p = 2
```

M =

	1	2	3	4
1	6	27	25	15
2	6	23	28	13
3	18	64	36	22
4	8	44	35	29
5	11	30	15	31
6	34	75	44	64
7	28	26	42	30
8	71	124	54	64
9	43	54	34	56
10	33	30	29	20
11	20	14	39	21

X =

	1	2
1	6	27
2	6	23
3	18	64
4	8	44
5	11	30
6	34	75
7	28	26
8	71	124
9	43	54
10	33	30
11	20	14

Y =

	1	2
1	25	15
2	28	13
3	36	22
4	35	29
5	15	31
6	44	64
7	42	30
8	54	64
9	34	56
10	29	20
11	39	21

Calculate Difference Matrix:

$$D := X - Y$$

Summary Statistics:

$$i := 1..n \quad j := 1..p$$

$$I := \text{identity}(n) \quad 1_{n_i} := 1$$

Mean Vectors:

$$X_{\text{bar}} := \frac{1}{n} \cdot X^T \cdot 1_n \quad X_{\text{bar}} = \begin{pmatrix} 25.272727 \\ 46.454545 \end{pmatrix}$$

$$Y_{\text{bar}} := \frac{1}{n} \cdot Y^T \cdot 1_n \quad Y_{\text{bar}} = \begin{pmatrix} 34.636364 \\ 33.181818 \end{pmatrix}$$

$$D_{\text{bar}} := \frac{1}{n} \cdot D^T \cdot 1_n \quad D_{\text{bar}} = \begin{pmatrix} -9.363636 \\ 13.272727 \end{pmatrix}$$

D =

	1	2
1	-19	12
2	-22	10
3	-18	42
4	-27	15
5	-4	-1
6	-10	11
7	-14	-4
8	17	60
9	9	-2
10	4	10
11	-19	-7

Covariance Matrices:

$$S_X := \frac{1}{n-1} \cdot X^T \cdot \left(I - \frac{1}{n} \cdot 1_n \cdot 1_n^T \right) \cdot X \quad S_X = \begin{pmatrix} 387.41818 & 489.36364 \\ 489.36364 & 1014.07273 \end{pmatrix}$$

$$S_Y := \frac{1}{n-1} \cdot Y^T \cdot \left(I - \frac{1}{n} \cdot 1_n \cdot 1_n^T \right) \cdot Y \quad S_Y = \begin{pmatrix} 109.25455 & 120.37273 \\ 120.37273 & 363.76364 \end{pmatrix}$$

$$S_D := \frac{1}{n-1} \cdot D^T \cdot \left(I - \frac{1}{n} \cdot 1_n \cdot 1_n^T \right) \cdot D \quad S_D = \begin{pmatrix} 199.25455 & 88.30909 \\ 88.30909 & 418.61818 \end{pmatrix}$$

Assumptions:

- Observed values $D_1, D_2, D_3, \dots, D_n$ are a random sample from $\sim N_p(\mu, \Sigma_D)$.
- Covariance Matrix Σ of the population is *unknown* and estimated by S_D .

Hypotheses:

Specify a test vector μ_0 :

$$\delta_0 := \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$H_0: \delta = \delta_0$ $< \delta_0 (0,0)$ or any value of difference between X & Y is a specified value for δ

$H_1: \delta \neq \delta_0$ $< \text{Note: Unlike univariate t-tests, All Hotelling } T^2 \text{ tests are two-sided since the } T^2 \text{ distribution, or equivalent F distribution, is asymmetric.}$

\wedge Note: this test is reasonably robust for deviations from $\sim N_p(\mu, \Sigma_D)$ with sufficient sample size

Hotelling's Test Statistic:

$$T_{sq} := n \cdot (D_{bar} - \delta_0)^T \cdot (S_D)^{-1} \cdot (D_{bar} - \delta_0) \quad T_{sq} = (13.639312)$$

$\wedge T^2$ is the normalized distance between vectors D_{bar} and δ_0

Likelihood Ratio Test Statistic:

Maximum likelihood estimates of μ & Σ (JW Result 4.11 p. 171) and Σ_0 :

$$\mu_{hat} := X_{bar}$$

$$\Sigma_{hat} := \frac{(n-1)}{n} \cdot S_D \quad |\Sigma_{hat}| = 6.249 \times 10^4$$

$$\Sigma_{0,hat} := \sum_i \left[\begin{pmatrix} X^T \\ \hat{\mu} \end{pmatrix} - \delta_0 \right] \cdot \left[\begin{pmatrix} X^T \\ \hat{\mu} \end{pmatrix} - \delta_0 \right]^T \quad |\Sigma_{0,hat}| = 5.216 \times 10^7$$

Likelihood Ratio & Wilks' lambda (Λ) (JW Eq. 5-13 p. 217):

$$\Lambda := \left(\frac{|n \cdot \Sigma_{hat}|}{|\Sigma_{0,hat}|} \right)^{\frac{n}{2}} \quad \Lambda = 0.00002438 \quad \Lambda^{\frac{2}{n}} = 0.144974$$

$< \Lambda = \text{Wilks' lambda - value is 1.0 when } X_{bar} = \mu_0 \text{ but decreases as } \mu_0 \text{ increases in distance from the sample mean.}$

Converting to Hotellings T2:

$$\left[1 + \frac{T_{sq}}{(n-1)} \right]^{-1} = (0.423024) \quad < \text{Equivalent value in terms of Hotellings } T^2 \text{ (jw Result 5.1 p. 218)}$$

$$\left[\Lambda^{-\left(\frac{2}{n}\right)} - 1 \right] \cdot (n-1) = 58.978$$

$$T_{sq} = (13.639)$$

$< \text{Solving for } T^2 \text{ (same as } T^2 \text{ above)}$

Sampling Distribution:

If Assumptions hold and H_0 is true, then $T_{sq} \sim T^2_{(n-1)} = [(n-1)p/(n-p)] F_{(p,n-p)}$

Critical Value of the Test:

$\alpha := 0.05$ < Probability of Type I error must be explicitly set

$$C := \frac{(n-1) \cdot p}{(n-p)} \cdot qF(1-\alpha, p, n-p) \quad C = 9.459 \quad < \text{JW eq. 5-6 p. 212}$$

NOTE: qF(1- α) is used in function qF.

Decision Rule:

Reject H_0 if $T_{sq} > C$

$$T_{sq} = (13.639) \quad C = 9.4589$$

Decision := if($T_{sq_1} > C, 1, 0$)

$$\text{Decision} = 1 \quad < 0 = \text{Do not reject } H_0 \\ 1 = \text{Reject } H_0$$

^ Therefore DO NOT Reject H_0

Probability Value:

$$P := 1 - pF\left[T_{sq} \cdot \frac{(n-p)}{(n-1) \cdot p}, p, n-p\right] \quad P = (0.02082779)$$

The Multivariate Confidence Ellipsoid:

$\lambda := \text{reverse}(\text{sort}(\text{eigenvals}(S_D)))$

$\varepsilon^{\langle j \rangle} := \text{eigenvec}(S_D, \lambda_j)$

$$\lambda = \begin{pmatrix} 449.75041 \\ 168.12231 \end{pmatrix} \quad \varepsilon = \begin{pmatrix} 0.33248 & -0.94311 \\ 0.94311 & 0.33248 \end{pmatrix}$$

< Coordinates of each column vector of ε gives the directions of confidence ellipsoid

$$CT := \sqrt{\frac{p \cdot (n-1)}{n \cdot (n-p)} \cdot qF(1-\alpha, p, n-p)} \quad CT = 0.9273068$$

< CT gives the boundary for the confidence ellipsoid for μ - see JW eq. 5-18 p. 221

$i := 1..p$

$$L_i := CT \cdot \sqrt{\lambda_i}$$

Multivariate confidence ellipsoid (JW Eq. 5-18 p. 221):

$$D_{\text{bar}} = \begin{pmatrix} -9.364 \\ 13.273 \end{pmatrix} \quad < \text{Center of ellipsoid} \quad L = \begin{pmatrix} 19.66569 \\ 12.02364 \end{pmatrix}$$

< L are half-lengths of the axes of the confidence ellipsoid for m in the directions of ε

Simultaneous T² Confidence Intervals:

$$CI_{\text{lower}_i} := D_{\text{bar}_i} - \sqrt{C} \cdot \sqrt{\frac{S_{D_{i,i}}}{n}} \quad CI_{\text{upper}_i} := D_{\text{bar}_i} + \sqrt{C} \cdot \sqrt{\frac{S_{D_{i,i}}}{n}} \quad < \text{jw eq. 5-24 p. 225 (slightly modified)}$$

CI := augment(CI_{lower} , CI_{upper})

Simultaneous Confidence Intervals:

$$D_{\text{bar}} = \begin{pmatrix} -9.36364 \\ 13.27273 \end{pmatrix} \quad < \text{Mean values} \quad CI = \begin{pmatrix} -22.45327 & 3.726 \\ -5.70012 & 32.24557 \end{pmatrix} \quad < T^2 \text{ confidence intervals}$$

Univariate t Confidence Intervals:

$$ct := qt\left(1 - \frac{\alpha}{2}, n - 1\right) \quad ct = 2.228139$$

< Critical value ct based on t distribution without correction

$$ci_{lower_i} := D_{bar_i} - ct \cdot \sqrt{\frac{SD_{i,i}}{n}} \quad ci_{upper_i} := D_{bar_i} + ct \cdot \sqrt{\frac{SD_{i,i}}{n}}$$

< jw eq. 5-29 p. 232

$$ci := augment(ci_{lower}, ci_{upper})$$

Univariate t Intervals:

$$D_{bar} = \begin{pmatrix} -9.36364 \\ 13.27273 \end{pmatrix} \quad < \text{Mean values} \quad ci = \begin{pmatrix} -18.84673 & 0.11946 \\ -0.4726 & 27.01805 \end{pmatrix}$$

< Univariate t confidence intervals

Bonferroni Simultaneous Confidence Intervals:

$$cb := qt\left(1 - \frac{\alpha}{2 \cdot p}, n - 1\right) \quad cb = 2.633767$$

< Critical value cb based on t distribution with Bonferroni correction factor p.

$$ci_{lower_i} := D_{bar_i} - cb \cdot \sqrt{\frac{SD_{i,i}}{n}} \quad ci_{upper_i} := D_{bar_i} + cb \cdot \sqrt{\frac{SD_{i,i}}{n}}$$

< jw eq. 5-29 p. 232

$$ci := augment(ci_{lower}, ci_{upper})$$

Bonferroni Intervals:

$$D_{bar} = \begin{pmatrix} -9.36364 \\ 13.27273 \end{pmatrix} \quad < \text{Mean values} \quad ci = \begin{pmatrix} -20.573107 & 1.845835 \\ -2.974903 & 29.520358 \end{pmatrix}$$

< Bonferroni confidence intervals

Prototype in R:

#HOTELLING'S T2 TEST FOR PAIRED DESIGN:

#LOAD DATA & HYPOTHESIS:

M=read.table("/DATA/Multivariate/JW6-1-longform.txt")

M #DATA TABLE

#SEPARATING GROUPS:

attach(M)

X=cbind(BOD=BOD[Lab==1],SS=SS[Lab==1])

X

Y=cbind(BOD=BOD[Lab==2],SS=SS[Lab==2])

Y

detach(M)

> M #DATA TABLE

	BOD	SS	Lab
1	6	27	1
2	6	23	1
3	18	64	1
4	8	44	1
5	11	30	1
6	34	75	1
7	28	26	1
8	71	124	1
9	43	54	1
10	33	30	1
11	20	14	1
12	25	15	2
13	28	13	2
14	36	22	2
15	35	29	2
16	15	31	2
17	44	64	2
18	42	30	2
19	54	64	2
20	34	56	2
21	29	20	2
22	39	21	2

Note the data structure in R is the so-called "long form". This is by far the most common way to handle data comprising multiple groups. The "dummy variable" (or factor) Lab is used to indicate group membership, with each variable comprising only a single long column.

#FUNCTION FOR HOTELLING'S T2 TEST FOR PAIRED DESIGN:

```
#X = dataset 1
#Y = dataset 2
#mu0 = hypothesis vector
#alpha = alpha of the test
```

```
PD.Hotelling.T2 <- function(X,Y,mu0,alpha)
```

```
... body of function in R script ...
```

```
PD.Hotelling.T2(X,Y,mu0=c(0,0),alpha=0.05)
```

```
RES=PD.Hotelling.T2(X,Y,mu0=c(0,0),alpha=0.05)
```

```
RES
```

```
> X
      BOD  SS
[1,]    6  27
[2,]    6  23
[3,]   18  64
[4,]    8  44
[5,]   11  30
[6,]   34  75
[7,]   28  26
[8,]   71 124
[9,]   43  54
[10,]  33  30
[11,]  20  14

> Y
      BOD  SS
[1,]   25 15
[2,]   28 13
[3,]   36 22
[4,]   35 29
[5,]   15 31
[6,]   44 64
[7,]   42 30
[8,]   54 64
[9,]   34 56
[10,]  29 20
[11,]  39 21
```

```
> RES=PD.Hotelling.T2(X,Y,mu0=c(0,0),alpha=0.05)
```

```
Paired Design Hotelling's T2
```

```
Hypothesis Vector:      ( 0 0 )
Difference Vector       ( -9.363636 13.27273 )
Discriminant vector   : ( -0.0673414 0.04591197 )
T2 Ellipsoid half lengths: ( 19.66569 12.02364 )
Hotelling's T2 Statistic: 13.63931
Equivalent F Statistic : 6.13769
F degrees of freedom:  ( 2 9 )
Wilks's Lambda:       0.008810659
alpha:                 0.05
Critical Value:        9.458877
Probability:           0.02082779
```

```
Confidence Intervals: T2 - Bonferroni - t - Mean - t - Bonferroni - T2
```

	T2.lower	B.lower	t.lower	Mean	t.upper	B.upper	T2.upper
BOD	-22.453272	-20.573107	-18.8467298	-9.363636	0.119457	1.845835	3.72600
SS	-5.700119	-2.974903	-0.4725958	13.272727	27.018050	29.520358	32.24557

#FUNCTION PLOT CI:

```
#R = return from function OS.Hotelling.T2 as a list
```

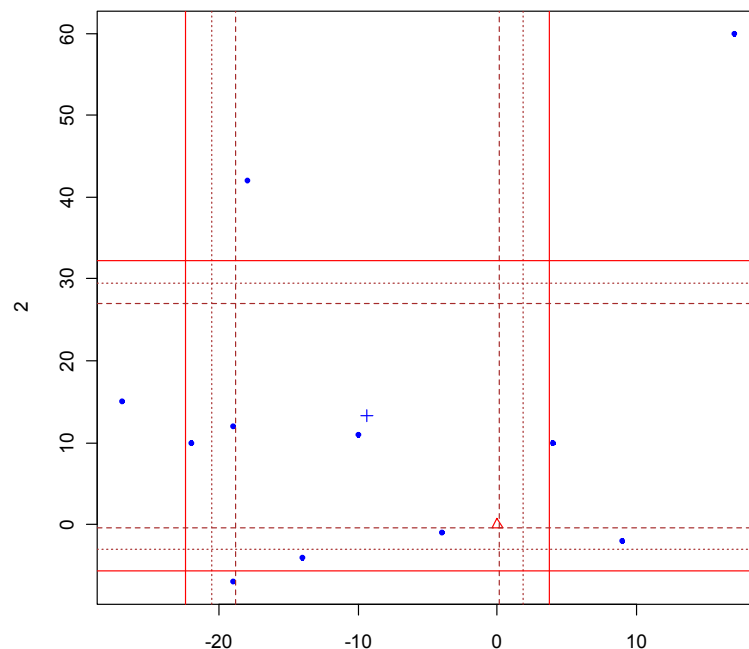
```
#X = variable to plot along the X axis
```

```
#Y = variable to plot along the Y axis
```

```
plot.CI <- function(R,X,Y)
```

```
... body of function in R script ...
```

```
plot.CI(RES,1,2)
```



1

Simultaneous T² intervals (solid red) are the widest, univariate t intervals (dashed) are the narrowest, with Bonferroni intervals (dotted) intermediate. Blue cross is sample mean vector. red triangle is hypothesis vector: