

14] Let $w_{ij} = k(x_i, x_j) = w_{ji}$, $i \neq j$.

HD 69

Let $H_n = \sum_{i < j} w_{ij}$. Then

$U_n = \frac{1}{\binom{n(n-1)}{2}}$ H_n is a U -statistic.
Sample mean of the w_{ij}

Note that $T_n = \frac{2 H_n}{n(n-1)} = U_n$ with $w_{ij} = \underline{x}_i^T \underline{x}_j$.

Aug 2023!

15] Theorem: Let $\underline{x}_1, \dots, \underline{x}_n$ be iid and $w_{ij} = \underline{x}_i^T \underline{x}_j$ for $i \neq j$. Let $E(\underline{x}_i) = \underline{\mu}$ and $V(w_{ij}) = \sigma_w^2$. Let

$\Theta = \text{cov}(w_{ij}, w_{id})$ where $i \neq j, i < j$ and $i < d$. Then

$$a) V(T_n) = \frac{2\sigma_w^2}{n(n-1)} + \frac{4(n-2)}{n(n-1)} \Theta$$

usually this term dominates

with $\Theta \geq 0$.

b) If $H_0 \underline{\mu} = \underline{0}$ is true, then

$$V(T_n) = \frac{2\sigma_w^2}{n(n-1)}$$

since $\Theta = 0$

$$16) \hat{\sigma}_m^2 = \frac{1}{m-1} \sum_{i=1}^m (w_i - \bar{w})^2 \xrightarrow{P} \sigma_w^2$$

69.5

where the w_i are as in 13].

$$\hat{\sigma}_w^2 = \frac{1}{n(n-1)} \sum_{i \neq j} (w_{ij} - T_n)^2$$

is nearly the sample variance of the w_{ij} , which are not ind.

from U statistics theory

17) Partial proof of 15].

$$V(H_n) = \text{cov}(H_n, H_n) =$$

$$\text{cov} \left(\sum_{i < j} w_{ij}, \sum_{k < d} w_{kd} \right) =$$

$$\sum_{i < j} \sum_{k < d} \text{cov}(w_{ij}, w_{kd}).$$

(*)

The covariances are of 3 types:

I) $(ij) = (kd)$ with $i < j$:

$$\text{cov}(w_{ij}, w_{kd}) = V(w_{ij}) = \sigma_w^2$$

II) i, j, k, d are distinct with $i < j, k < d$.

Then $w_{ij} \perp w_{kd}$ and $\text{cov}(w_{ij}, w_{kd}) = 0$.

III) Exactly three of the four subscripts are distinct which have $cov(w_{ij}, w_{kl}) = \theta$ for $j \neq d, i < j$, and $i < d$ or $i \neq k, i < j, k < d$ which have $cov(w_{ij}, w_{kl}) = \theta$ since $w_{ij} = w_{ji}$.

The number of ways to get 3 distinct subscripts is $a - b - c$

Done in U-statistics theory

$$= \binom{n}{2} - \underbrace{\binom{n}{2} \binom{n-2}{2}}_{b} - \underbrace{\binom{n}{2}}_c = n(n-1)(n-2)$$

$a = \# \text{ terms in } \binom{n}{2}$
 $b = \# \text{ terms where } i, j, k, d \text{ are distinct with } i < j \text{ and } k < d$
 $c = \# \text{ terms with } (ij) = (kd) \text{ and } i < j.$

Thus $V(H_n) = 0.5 n(n-1) \sigma_w^2 + n(n-1)(n-2) \theta$

and $V(T_n) = \frac{4}{[n(n-1)]^2} V(H_n) = \frac{2 \sigma_w^2}{n(n-1)} + \frac{4(n-2)\theta}{n(n-1)}$

$\theta \geq 0$ also comes from U-statistics theory.

18] Let $s_0 = \sqrt{\frac{2 \hat{\sigma}_w^2}{n(n-1)}}$ or $\sqrt{\frac{2 s_m^2}{n(n-1)}}$ from 16], 70.5

These formulas are simpler than the ones in the published literature.

When H_0 is not true,
 suppose $s_n^2 \xrightarrow{P} v(T_n)$.

$$\frac{T_n}{s_0} \approx \frac{T_n}{s_n} \quad \text{under } H_0,$$

but $\frac{|T_n|}{s_0} > \frac{|T_n|}{s_n}$ since $\theta > 0$ when

H_0 is false. Thus $\frac{|T_n|}{s_0}$ has

more power when H_0 is false, where

power $\equiv P(\text{reject } H_0)$, (want power near level (eg. 0.05) when H_0 is true, power near 1 if H_0 is false)

19] HD tests for $H_0 \underline{\mu}_1 = \underline{\mu}_2$ vs $H_1 \underline{\mu}_1 \neq \underline{\mu}_2$.

want to estimate $\|\underline{\mu}_1 - \underline{\mu}_2\|^2 =$

$$(\underline{\mu}_1 - \underline{\mu}_2)^T (\underline{\mu}_1 - \underline{\mu}_2) = \underline{\mu}_1^T \underline{\mu}_1 + \underline{\mu}_2^T \underline{\mu}_2 - 2 \underline{\mu}_1^T \underline{\mu}_2.$$

Suppose we have 2 independent
(iid) random samples

HD 71

$$\underline{x}_{1,1}, \dots, \underline{x}_{1,n_1}$$

$$\underline{x}_{2,1}, \dots, \underline{x}_{2,n_2}$$

$$\text{Let } \underline{a} = \sum_{i=1}^{n_1} \underline{x}_{1,i}, \quad \underline{X}_1 = (\underline{x}_{1,i}) = \begin{pmatrix} \underline{x}_{1,1}^T \\ \vdots \\ \underline{x}_{1,n_1}^T \end{pmatrix}$$

$$\underline{c} = (\underline{x}_{1,1}^T, \dots, \underline{x}_{1,n_1}^T)^T$$

$$\underline{b} = \sum_{i=1}^{n_2} \underline{x}_{2,i}, \quad \underline{X}_2 = (\underline{x}_{2,i}) = \begin{pmatrix} \underline{x}_{2,1}^T \\ \vdots \\ \underline{x}_{2,n_2}^T \end{pmatrix}$$

$$\underline{d} = (\underline{x}_{2,1}^T, \dots, \underline{x}_{2,n_2}^T)^T$$

$$\text{Let } T_n = \frac{1}{n_1(n_1-1)} [\underline{a}^T \underline{a} - \underline{c}^T \underline{c}] + \frac{1}{n_2(n_2-1)} [\underline{b}^T \underline{b} - \underline{d}^T \underline{d}]$$

$$- \frac{2 \underline{a}^T \underline{b}}{n_1 n_2}$$

Under $H_0: \underline{\mu}_1 = \underline{\mu}_2$

and $s_n^2 \xrightarrow{P} V(T_n)$

$$\frac{T_n}{s_n} \xrightarrow{D} N(0, 1)$$

PhD topic:

71.5

13-193 try to get better S_n for 193.

20] "m out of n" bootstrap with $m_i = \lfloor \frac{2n_i}{3} \rfloor$ simulated well.

21] For the 1-sample problem,

$$T_m^* = \frac{1}{m(m-1)} \sum_{k \neq d}^m \sum_{i \neq d}^m x_{-k}^T x_{-i}$$
 where

x_{-1}, \dots, x_{-m} are the cases selected without replacement. $n(n-1)$ terms in the sum

$$\text{Hence } T_m^* = \frac{1}{m(m-1)} \sum_{i \neq j}^n w_i w_j \underbrace{x_i^T x_j}_{\text{constant w/rt bootstrap dist.}}$$

where the indicator $w_k = \begin{cases} 1 & x_k \text{ in the sample} \\ 0 & \text{else} \end{cases}$

The w_k are identically distributed but dependent since $\sum_{k=1}^n w_k = m$.

Since the w_k are binary $\in \{0,1\}$, HD72

$$w_k^a = w_k \text{ for } a > 0.$$

A product of indicators is an indicator
 $= 1$ iff all of the indicators $= 1$.

$$\text{Hence } w_{i_1}^{a_1} \dots w_{i_k}^{a_k} = w_{i_1} \dots w_{i_k}$$

$= w_{i_1, i_2, \dots, i_k}$ is the indicator for

whether x_{i_1}, \dots, x_{i_k} are all in the
sample of size m . Here i_1, \dots, i_k are
 k distinct indices.

$$\text{Thus } E[w_{i_1}^{a_1} \dots w_{i_k}^{a_k}] = E(w_{i_1} \dots w_{i_k})$$

$$= \frac{m}{n} \frac{m-1}{n-1} \dots \frac{m-k+1}{n-k+1} \quad \left. \vphantom{\frac{m}{n}} \right\} \text{ slots}$$

permutations

Note that $m(m-1)\dots(m-k+1)$ ordered k tuples
were selected from $n(n-1)\dots(n-k+1)$ possible
ordered k tuples (sequences of length k).

22) Thus $E[T_m^*] = \frac{1}{m(m-1)} \sum_{i \neq j} E(w_i w_j) x_i^T x_j$ (72.5)

$$= \frac{1}{m(m-1)} \frac{m(m-1)}{n(n-1)} \sum_{i \neq j} x_i^T x_j = T_n.$$

multiple testing

1) Suppose there is a single gene g and there are n measurements of "expression of gene g "

$$x_1^A, \dots, x_n^A$$

A = normal cells (control)

$$x_1^B, \dots, x_n^B$$

B = cancer cells (patient)

$$H_0 \mu_{Ag} = \mu_{Bg} \quad H_1 \mu_{Ag} \neq \mu_{Bg}$$

2) Suppose there is a test statistic for gene g that produces (an estimated) p-value \hat{p}_g (eg reject H_0 at level α if $\hat{p}_g < \alpha$).

3) With DNA microarrays and next generation sequencing techniques,

Can measure the expression levels of m genes simultaneously, where m is in the thousands:

$$g \in \{1, \dots, m\}.$$

H_{0g} the mean expression levels of gene g in conditions A and B are the same

" different.

H_{1g} "

4) If $m > 20000$, $\alpha = .05$, and

H_{0g} is true for 20000 genes, then on average, there are $20000(.05) = 1000$ false positives (discoveries).

When H_{0g} is rejected further experiments need to be done in order to validate or not the "discovery" (that gene g is involved in the response to change of "environment" between A and B). These experiments cost a lot of time and money. Want few false discoveries without losing too much in power.

5} Given m tests and m

p values $\hat{p}_1, \dots, \hat{p}_m$, let

$\hat{R} = \hat{R}(\hat{p}_1, \dots, \hat{p}_m)$ give a set of indices i , for hypotheses

H_{0i} , that are rejected,

$\hat{R} = \emptyset$ is possible.

6} Let $I_0 = \{i \in \{1, \dots, m\} : H_{0i} \text{ is true}\}$.

we call false positive the

indices $i \in \hat{R} \cap I_0$ and true

positive

the indices $i \in \hat{R} \cap I_0^c = \hat{R} \setminus I_0$.

Let $FP = \text{card}(\hat{R} \cap I_0) = \#$ of false positives

and $TP = \text{card}(\hat{R} \cap I_0^c) = \#$ of true positives,

Here $I_0^c = \{i \in I : i \notin I_0\}$.

7} The Bonferroni correction

uses $\hat{R}_B = \{i : \hat{p}_i \leq \frac{\alpha}{m}\}$.

8) Some assumptions on the p-value \hat{p}_i are needed, HD 74

such as $\sup_{\theta \in \Theta_0} P_{\theta}(\hat{p}_i \leq u) \leq u \quad \forall u \in [0, 1]$,

9) If 8) holds, then for \hat{R}_B \leftarrow

$$P(FP > 0) \leq \alpha.$$

Thus the existence of false positives $\leq \alpha$, but TP is often small.

10) The false discovery proportion

$$FDP = \frac{FP}{FP + TP} \quad \text{with } \frac{0}{0} = 0.$$

The false discovery rate

$$FDR = E \left[\frac{FP}{FP + TP} \mathbb{I}(FP + TP \geq 1) \right]$$

and is widely used in biostatistics.

11) To make FP as small as possible and TP as large as possible, the

$i \in \hat{R}$ correspond to the smallest p-values.

Need to determine how many of the smallest

\hat{p}_i should be in \hat{R} .

74.5

want card(\hat{R} / \hat{K} "large" with $FDR \leq \alpha$,

so $\hat{p}_{(1)}, \dots, \hat{p}_{(\hat{K})} \in \hat{R}$ for some \hat{K} ,

where $\hat{R} = \emptyset$ is possible ($\hat{K} = 0$).

12) Consider the line

$$\hat{p} = \gamma = 0 + \frac{\alpha}{m} k \quad \left(\text{where } 0 < \alpha < 1 \text{ is fixed} \right)$$

eg $\alpha = 0.05$

$$\text{Let } \hat{R}_{BH} = \emptyset \text{ if } \left\{ k \in I : \hat{p}_{(k)} \leq \frac{\alpha k}{m} \right\} = \emptyset.$$

$$\text{otherwise, let } \hat{K} = \max \left\{ k \in I : \hat{p}_{(k)} \leq \frac{\alpha k}{m} \right\}$$

$$\text{and } \hat{R}_{BH} = \left\{ i \in I : \hat{p}_i \leq \alpha \frac{\hat{K}}{m} \right\}$$

$$= \left\{ i_{(1)}, \dots, i_{(\hat{K})} \text{ corresponding to } \hat{p}_{(1)}, \dots, \hat{p}_{(\hat{K})} \right\}.$$

under fairly strong regularity conditions,

$$FDR \leq \alpha \text{ for } \hat{R}_{BH}.$$