1. **Problem**
   I have a sample $X_1, \ldots, X_{60}$. I generate one naive bootstrap sample $X_1^*, \ldots, X_{60}^*$.
   Let $N$ be the number of times the first observation will be copied in the bootstrap sample.
   Find the probability $\mathbb{P}(N = 2)$.

$$X_1 \rightrightarrows h \quad j^*$$

$$\begin{array}{cccc} \frac{1}{60} & \frac{1}{60} & \frac{59}{60} & \frac{59}{60} \\ 1 & 2 & 3 & \cdots \quad 58 \end{array}$$

$$X_1 \quad \nearrow \quad {}^1/60$$
$$\quad \searrow \quad 59/60$$

$$P(N = 2) = \frac{1}{60} \cdot \frac{1}{60} \cdot \left(\frac{59}{60}\right)^{58} \cdot C_{60}^2$$

$$N \sim \text{binom}(n, p) = \text{Bin}\left(60, \frac{1}{60}\right)$$

$$P(N = k) = C_n^k \cdot p^n \cdot (1-p)^{n-k}$$

$$P(N \leq 1) = P(N = 0) + P(N = 1)$$

## 2. Problem

I have a sample $X_1, \ldots, X_{60}$. I generate one naive bootstrap sample $X_1^*, \ldots, X_{60}^*$.

Let $L$ be the number of initial observations missing in the bootstrap sample.

Find the expected value $\mathbb{E}(L)$.

$$\ell_i = \begin{cases} 1 & , \text{ missing} & 59/60 \\ 0 & , \text{ not missing} & 1/60 \end{cases}$$

$$\sum \ell_i = L$$

$$E(\ell_i) = p = \frac{59}{60} \qquad \text{Var}(\ell_i) = pq = \frac{59}{60} \cdot \frac{1}{60}$$

$$E(L) = \underbrace{1 \cdot \left(\frac{59}{60}\right)^{60}}_{X_1} + \underbrace{1 \cdot \left(\frac{59}{60}\right)^{60}}_{X_2} + \ldots + \underbrace{1 \cdot \left(\frac{59}{60}\right)^{60}}_{X_{60}}$$

$$= 60 \cdot \left(\frac{59}{60}\right)^{60}$$

3. **Problem**
We have data of an AB experiment: $\bar{X}_a = 5.1$, $\bar{X}_b = 5.6$, $n_a = 18$, $n_b = 12$, $\sum(X_i^a - \bar{X}_a)^2 = 840$, $\sum(X_i^b - \bar{X}_b)^2 = 820$.

Calculate the estimate of variance of $\bar{X}_a - \bar{X}_b$ for the Welch test.

$$t\text{-test} \qquad vs \qquad \text{Welch Test}$$

$$H_0: \quad \mu_a = \mu_b \qquad\qquad H_0: \quad \mu_a = \mu_b$$

$$\text{Ass.}: \quad \sigma_a = \sigma_b = \sigma_p \qquad \text{Ass.}: \quad \sigma_a \neq \sigma_b$$

$$\hat{v} \leq n-2 \qquad t = \frac{\bar{X}_a - \bar{X}_b - 0}{Se(\bar{X}_a - \bar{X}_b)} \sim t_{n-2} \qquad t = \frac{\bar{X}_a - \bar{X}_b - 0}{Se(\bar{X}_a - \bar{X}_b)} \sim t_{\hat{v}}$$

$$t = \frac{\hat{\theta} - \theta}{Se(\hat{\theta})}$$

$$Se(\bar{X}_a - \bar{X}_b) = \sqrt{\frac{S_p^2}{n_1} + \frac{S_p^2}{n_2}} \qquad Se(\bar{X}_a - \bar{X}_b) = \sqrt{\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}}$$

$$S_p^2 = \frac{TSS_1 + TSS_2}{n_1 + n_2 - 2} = \frac{(n_1-1)S_1^2 + (n_2-1)\cdot S_2^2}{n_1 + n_2 - 2}$$

$$\left[\begin{array}{l}\hat{\sigma}_i^2 = \frac{1}{n}\cdot \sum(X_i - \bar{X})^2 \\ S_1^2 = \frac{1}{n-1}\cdot \sum(X_i - \bar{X}_2)^2 = \end{array}\right.$$

$$= \frac{1}{n_1 - 1} TSS_1$$

$$S_2^2 = \frac{TSS_2}{n_2 - 1}$$

$$\hat{Var}(\bar{X}_a - \bar{X}_b) = \sqrt{\frac{S_a^2}{n_a} + \frac{S_b^2}{n_b}} = \sqrt{\frac{840/17}{18} + \frac{820/11}{12}}$$

## 4. Problem

Variables $X_1$, $X_2$, ..., $X_{20}$ are independent and normally distributed $\mathcal{N}(5;3)$.

Calculate the variance

$$\mathrm{Var}\left(\sum_{i=1}^{20}(X_i - \bar{X})^2\right).$$

$X_i \sim N(5,3) \quad 3 = \delta^2 \qquad \frac{\Sigma X_i}{n}$

$\mathrm{Cov}\left(X_i - \bar{X}, X_j - \bar{X}\right) \neq 0$

$\mathrm{Var}(TSS) - ?$

$\chi^2_k = \sum_{i=1}^{k} \zeta_i^2, \quad \zeta_i \sim N(0,1)$

$\Sigma \dfrac{(X_i - \bar{X})^2}{\delta^2} = \Sigma \zeta^2$

$S^2 = \dfrac{1}{n-1} \cdot TSS$

$\dfrac{(n-1) S^2}{\delta^2} \sim \chi^2_{n-1}$

$\hookrightarrow n \cdot R^2 \sim \chi^2$

$\hookrightarrow F = \dfrac{\mathrm{Var}(\cdot)}{\mathrm{Var}(\cdot)} \sim \dfrac{\chi^2_p}{\chi^2_q} = F_{p,q}$

$\mathrm{Var}\left(\dfrac{n-1}{\delta^2} S^2\right) = \mathrm{Var}(\chi^2_{n-1})$

$\dfrac{(n-1)^2}{\delta^4} \cdot \mathrm{Var}\left(\dfrac{TSS}{n-1}\right) = \mathrm{Var}(\chi^2_{n-1})$

$E(\zeta) = 0$

$E(\zeta^2) = \mathrm{Var}(\zeta)$

$\mathrm{Var}(\chi^2_1) = \mathrm{Var}(\zeta^2) =$

$E(\zeta^4) - E^2(\zeta^2) = \mu_4 - \mu_2^2 =$

$3 - 1 = 2$

$\mathrm{Var}(\chi_k) = \mathrm{Var}\left(\sum_{i=1}^{k} \zeta_i^2\right) = \sum_{i=1}^{k} \mathrm{Var}(\zeta_i^2) =$

$= k \cdot 2$

$\dfrac{\mathrm{Var}(TSS)}{\delta^4} = (n-1) \cdot 2 \qquad \mathrm{Var}(TSS) = 2 \cdot (n-1) \cdot \delta^4 =$

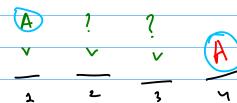$= 2 \cdot 19 \cdot \left(\sqrt{3}\right)^4$

## 5. Problem

I have results of two runners A and B: 2 results for A and 3 results for B. The running time for both runners are continuously distributed and their distribution are equal.

What is the probability that the maximal rank of running times of the runner A will be equal to 4?

$$X^A \sim N(\mu_a, \sigma_a^2)$$

$$\text{rank}(X_3^B) = 5$$

(\*)
$$\begin{cases} A: & 2 \quad 4 \\ B: & 3 \quad 5 \quad \boxed{6} \end{cases}$$

$$\boxed{A} \quad ? \quad ? \quad \boxed{A} \quad B$$
$$\underset{1}{\vee} \quad \underset{2}{\vee} \quad \underset{3}{\vee} \quad \overline{4} \quad \overline{5}$$

$$C_1^3 \quad : \quad ABB \quad BAB \quad BBA$$

$$P(\max(\text{rank } A) = 4) = \frac{C_1^3}{C_2^5}$$

$$P(\max(\text{rank } A) \leq 4) \quad \nearrow \quad \text{rank} = 1 \quad ?? \, ABB$$
$$A$$

$$\text{rank} = 2 \quad AABBB$$

$$P(\text{rank} = 1) = 0$$

## 6. Problem

I have results of two runners A and B: 2 results for A and 3 results for B. The running times for both runners are continuously distributed.

Consider the Mann-Whitney test statistic $U_a$ that positively depends on the rank sum for the runner A.

Find the expected value $\mathbb{E}(A)$ under the null-hypothesis of equal running time distributions.

T/W $\qquad$ $H_0:\ \mu_a = \mu_b$

$A:\quad 2\quad 4$ $\qquad$ MW $\qquad H_0: E(\text{rank}_a) = E(\text{rank}_b)$

$B:\quad\quad 3\ 5\ (6)$ $\qquad \xrightarrow{\ \ \vec{c}\ \ }$ $\quad \hookrightarrow H_0: f_a(x) = f_b(x+c),\ c=0$

$\qquad\qquad\qquad\qquad\qquad\qquad \hookrightarrow H_0: P(X \leqslant y) = 0.5$

$\qquad\qquad\qquad\qquad\qquad\qquad \hookrightarrow H_0: \mu_a = \mu_b$

|   | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|
| A | 0 |   | 1 |   |   |
| B |   | 1 |   | 2 | 2 |

$\Sigma \Rightarrow U_a = 1$

$\Sigma \Rightarrow U_b = 5$

$$U = \min\{U_a, U_b\} = 1$$

$$U \stackrel{d}{\sim} N\left(\frac{n_1 n_2}{2};\ \frac{n_1 n_2 (n_1 + n_2 - 1)}{12}\right)$$

$U_{ob} = 1$ $\qquad$ vs $\qquad U_{crit}$

$$\frac{U_{obs} - E(U)}{sd(U)} \sim N(0,1)$$

$$\frac{1 - \frac{2 \cdot 3}{2}}{\sqrt{\frac{2 \cdot 3 \cdot 4}{12}}} = \frac{-2}{\sqrt{2}} \approx -1.4 \qquad \left(|-1.4| < 1.96\right.$$

$$\Rightarrow H_0 \text{ is not rej.}$$

## 7. Problem

Consider three variables: target variable $y_i$, predictor $x_i$ and indicator of treatment $z_i \in \{0,1\}$.

The treatment $z_i$ was designed to be independent of $x_i$, but in fact $x_i = f_i \cdot (1 + 0.08 z_i)$.

We suppose that $z_i$ are Bernoulli with $p = 0.4$, $f_i \sim \mathcal{N}(900; 9)$ and they are independent.

Find the probability limit

$$\text{plim} \frac{\sum (x_i - \bar{x})(z_i - \bar{z})}{n-1}.$$

$$\text{Cov}(X,X) = \text{Var}(X)$$

$$\plim_{n \to \infty} \frac{\bar{X}}{n} = E(X) = \mu$$

$$\plim_{n \to \infty} \sum (X_i - \bar{X})^2 = \infty$$

$$\plim_{n \to \infty} \frac{\sum (X_i - \bar{X})^2}{n} = \plim_{n \to \infty} \frac{\sum (X_i - \bar{X})}{n-1} = \text{Var}(X)$$

$$\plim_{n \to \infty} \widehat{\text{Cov}}(x,z) = \text{Cov}(X_i, Z_i) = \text{Cov}(f_i \cdot (1 + 0.08 z_i), z_i)$$

$$= \underbrace{\text{Cov}(f_i, z_i)}_{0} + 0.08 \cdot \text{Cov}(f_i z_i, z_i) =$$

$$z_i = \begin{cases} 1 & p \\ 0 & 1-p \end{cases}$$

$$z_i^2 = \begin{cases} 1^2 \\ 0^2 \end{cases}$$

$$= 0.08 \cdot \left( E(f_i z_i^2) - E(f_i z_i) \cdot E(z_i) \right) =$$

$$= 0.08 \cdot E(f_i z_i) \left( 1 - E(z_i) \right) =$$

$$= 0.08 \cdot 900 \cdot 0.4 \cdot (1 - 0.4)$$

## 8. Problem

Consider three variables: target variable $y_i$, predictior $x_i$ and indicator of treatment $z_i \in \{0, 1\}$. The treatment $z_i$ was assigned independently of $x_i$, total $n = 100$.

The sample covariance matrix $C$ is provided. The order of variables is $y$, $x$ and $z$. For example the sample covariance of $x$ and $z$ is 0.26.

$$C = \begin{matrix} & y & x & z \\ y & \begin{pmatrix} 7.09 & 2.64 & 0.75 \\ x & 2.64 & 1.17 & 0.26 \\ z & 0.75 & 0.26 & 0.25 \end{pmatrix} \end{matrix}$$

Consider CUPED with first regression given by $\hat{y}_i = \hat{\alpha}_1 + \hat{\alpha}_2 x_i$ with residuals $r_i = y_i - \hat{y}_i$.

What is the sample covariance of $r$ and $z$?

Cuped:

1) $Y_{cuped} = Y - \theta X$

$$\hat{\theta} = \frac{\hat{Cov}(X,Y)}{\widehat{Var}(X)} \qquad \leftrightarrow \quad \hat{\beta} \text{ in } Y|X$$

2) $Var(Y_{cuped}) = (1 - \rho_{x,y}^2) Var(Y) \leq Var(Y)$

$$\rho_{x,y}^2 = R_{y|x}^2$$

$$\hat{Cov}(r, z) = \hat{Cov}(y - \hat{y}, z) =$$

$$= \hat{Cov}(y - \hat{\theta}X, z) =$$

$$\hat{Cov}(y, z) - \hat{\theta} \, \hat{Cov}(x, z) =$$

$$\hat{Cov}(y, z) - \frac{\hat{Cov}(x, y)}{\widehat{Var}(X)} \hat{Cov}(x, z) =$$

$$0.75 - \frac{2.64}{1.17} \cdot 0.26$$