

# Continuous Buhmbox - Update 12/4

Alexandre Lamy - all2187@columbia.edu

December 4, 2017

## Notation

As a quick recall, remember that

$$S_{BUHMBOX} = \frac{\sum_{i < j} w_{ij} y_{ij}}{\sqrt{\sum_{i < j} w_{ij}^2}}$$

And,

$$Y = \alpha(R - R')$$

Where  $w_{ij}$  are weights and, in the standard BB case,  $R$  is the sample correlation matrix in cases and  $R'$  in controls, and in the continuous BB case,  $R$  is the weighted sample correlation matrix (using some weights  $\omega$ ) and  $R'$  is the unweighted sample correlation matrix (or equivalently the weighted sample correlation matrix with all equal weights).  $\alpha$  is simply a normalization factor to assure that all  $y_{ij}$  are  $\mathcal{N}(0, 1)$  in the independent case.

We denote  $r_\omega$  to be the weighted sample correlation between two independent random variables ( $X$  and  $Y$ ) calculated on  $N$  observations ( $x_i$  and  $y_i$ 's), each with weight  $\omega_i$  and with  $\sum_i \omega_i = 1$  (and  $\forall i, \omega_i \geq 0$ ). And  $r$  to be  $r_\omega$  where  $\omega$  just gives equal weights to all observations (equivalently the unweighted sample correlation), we sometimes also denote this as  $r_{1/N}$  for obvious reasons.

Then note that in the independent population case, each  $y_{ij}$  is distributed like an  $r_\omega - r$  for the continuous BB case.

Some extra notation:

- We let  $\bar{x} = \frac{1}{n} \sum_i x_i$  be the sample mean of observations.
- We let  $\hat{x}_\omega = \sum_i \omega_i x_i$  be the weighted sample mean of observations.
- We let  $\Omega_j = \sum_i \omega_i^j$  for a given weighting scheme  $\omega$ . Thus  $\Omega_1 = 1$ .

We also denote the sample mean of observation  $x_1, x_2, \dots, x_n$  as and the weighted sample mean in respect to some weighting  $\omega$  as .

## Theorems and empirically verified conjectures

I claim to have a proof for everything I label as a “Theorem”. I have empirically verified conjecture but do not yet have proofs.

**Theorem 1.** *Suppose that all  $w_{ij}$  are independent and prefixed weights, and assume all  $y_{ij}$  are  $\mathcal{N}(0,1)$ . Then  $S_{BUHMBOX}$  is distributed as a  $\mathcal{N}(0,1)$ .*

*Remark 1.* This theorem explains why proving that the  $y_{ij}$ ’s are  $\mathcal{N}(0,1)$  is important. However, the strong (and possibly wrong) assumption of the independence of the  $w_{ij}$ ’s might be the cause of the main current challenge (see remaining problems/challenges section). The rest of the theorems involve trying to prove that the  $y_{ij}$ ’s (or equivalently  $r_\omega$  for arbitrary  $\omega$ ’s) are distributed as  $\mathcal{N}(0,1)$  under independence.

**Theorem 2.**  $\mathbb{E}[r_\omega] = 0$  for any valid  $\omega$  (by valid we mean  $\sum_i \omega_i = 1$  and  $\forall i, \omega_i \geq 0$ ).

*Remark 2.* This along with the linearity of expectation and the fact that  $r$  is just a special case of  $r_\omega$ , proves that  $\mathbb{E}[y_{ij}] = 0$  in both the standard and continuous versions of Buhmbox.

**Theorem 3.**  $\text{Var}[r] = \mathbb{E}[r^2] \approx \frac{1}{N}$  for large  $N$ .

*Remark 3.* This comes trivially from the fact that, for normal variables  $X$  and  $Y$ , the sample correlation coefficient has a variance of approximately

$$\frac{1 - \rho^2}{N - 2}$$

which is approximately  $\frac{1}{N}$  in our case. That expression from the variance comes from long papers that use the Fischer Transformation. However, we successfully proved this from first principles, without even relying on the assumption of normality. When doing this we proved that the variance was exactly  $\frac{1}{N-1}$

**Conjecture 1.**  $\text{Var}[r_\omega] = \mathbb{E}[r_\omega^2] \approx \sum_i \omega_i^2$

*Remark 4.* Note that this is a generalization of theorem 3. However, although being very close, we have not been able to prove this result using similar techniques (I get stuck with a really long expression, I can’t simplify). This may be due to the above result being an approximation rather than an exact result. However, it is empirically an **excellent** approximation.

**Conjecture 2.**  $\text{Var}[r_\alpha - r_\beta] = \mathbb{E}[(r_\alpha - r_\beta)^2] \approx \sum_i (\alpha_i - \beta_i)^2$ .

Or, equivalently  $\mathbb{E}[r_\alpha r_\beta] \approx \sum_i \alpha_i \beta_i$ .

*Remark 5.* First note that this is again a generalization of both theorem 3 and conjecture 1. Furthermore it was verified to be empirically extremely accurate and matches exactly with the formulation of standard Buhmbox. Indeed note that the standard Buhmbox just calculates  $y_{ij}$  as  $r_\alpha - r_\beta$  where  $\alpha$  gives weights of  $\frac{1}{N}$  for all cases and weights of 0 for all controls, and  $\beta$  gives weights of  $\frac{1}{N'}$  to all controls and weights of 0 to all cases. In which case this result gives us that:

$$\text{Var}[y_{ij}] = \mathbb{E}[(r_\alpha - r_\beta)^2] = \sum_i (\alpha_i - \beta_i)^2 = \sum_{\text{cases}} \left( \frac{1}{N} - 0 \right)^2 + \sum_{\text{controls}} \left( 0 - \frac{1}{N'} \right)^2 = N \frac{1}{N^2} + N' \frac{1}{N'^2} = \frac{1}{N} + \frac{1}{N'}$$

Which exactly matches their normalization factor of  $\sqrt{\frac{NN'}{N+N'}} = \sqrt{\frac{1}{\frac{1}{N} + \frac{1}{N'}}}$ .

As another quick note, equivalence of the two results simply comes from the fact that

$$\text{Var}[r_\alpha - r_\beta] = \text{Var}[r_\alpha] + \text{Var}[r_\beta] - 2\text{Cov}[r_\alpha, r_\beta] = \text{Var}[r_\alpha] + \text{Var}[r_\beta] - 2\mathbb{E}[r_\alpha r_\beta]$$

And a teeny bit of algebra.

Proof wise, the first result seems almost impossible to prove directly. The second seems easier, but proving conjecture 1 would be a first step. There is also an added complexity in that the denominator has a square root that I can't get rid of, and taking the expectation of a square root seems particularly difficult.

## Lemmas

A bunch of lemmas that we will rely on in our proofs:

**Lemma 1.**  $\mathbb{E}[\hat{x}_\omega] = \mathbb{E}[X]$

*Proof.*

$$\begin{aligned}\mathbb{E}[\hat{x}_\omega] &= \mathbb{E}\left[\sum_{k \in I} \omega_k x_k\right] \\ &= \sum_{k \in I} \omega_k \mathbb{E}[x_k] \\ &= \sum_{k \in I} \omega_k \mathbb{E}[X] \\ &= \mathbb{E}[X] \sum_{k \in I} \omega_k \\ &= \mathbb{E}[X] \Omega_1 \\ &= \mathbb{E}[X]\end{aligned}$$

□

**Lemma 2.**  $\mathbb{E}[\bar{x}^2] = \frac{\text{Var}[X]}{N} + \mathbb{E}[X]^2$

*Proof.*

$$\begin{aligned}\text{Var}[\bar{x}] &= \text{Var}\left[\frac{1}{N} \sum_i x_i\right] \\ &= \frac{1}{N^2} \text{Var}\left[\sum_i x_i\right] \\ &= \frac{1}{N^2} \sum_i \text{Var}[x_i] && \text{(by independence)} \\ &= \frac{1}{N^2} N \text{Var}[X] \\ &= \frac{1}{N} \text{Var}[X]\end{aligned}$$

And

$$\begin{aligned}\text{Var}[\bar{x}] &= \mathbb{E}[\bar{x}^2] - \mathbb{E}[\bar{x}]^2 \\ &= \mathbb{E}[\bar{x}^2] - \mathbb{E}[X]^2\end{aligned}$$

Hence  $\frac{1}{N} \text{Var}[X] = \mathbb{E}[\bar{x}^2] - \mathbb{E}[X]^2$  which gives the desired result.

□

**Lemma 3.**  $\mathbb{E}[x_i \bar{x}] = \frac{\text{Var}[X]}{N} + \mathbb{E}[X]^2$

*Proof.*

$$\begin{aligned}
 \mathbb{E}[x_i \bar{x}] &= \mathbb{E}\left[x_i \frac{1}{N} \sum_j x_j\right] \\
 &= \frac{1}{N} \mathbb{E}[x_i^2 + \sum_{j \neq i} x_i x_j] \\
 &= \frac{1}{N} (\mathbb{E}[x_i^2] + \sum_{j \neq i} \mathbb{E}[x_i x_j]) \\
 &= \frac{1}{N} (\mathbb{E}[X^2] + (N-1) \mathbb{E}[x_i] \mathbb{E}[x_j]) && \text{(By independence)} \\
 &= \frac{1}{N} (\mathbb{E}[X^2] + (N-1) \mathbb{E}[X]^2) \\
 &= \frac{1}{N} (\text{Var}[X] + N \mathbb{E}[X]^2) \\
 &= \frac{\text{Var}[X]}{N} + \mathbb{E}[X]^2
 \end{aligned}$$

□

**Lemma 4.**  $\mathbb{E}[(x_i - \bar{x})^2] = \frac{N-1}{N} \text{Var}[X]$

*Proof.*

$$\begin{aligned}
 \mathbb{E}[(x_i - \bar{x})^2] &= \mathbb{E}[x_i^2] + \mathbb{E}[\bar{x}^2] - 2 \mathbb{E}[x_i \bar{x}] \\
 &= \mathbb{E}[X^2] - \frac{\text{Var}[X]}{N} - \mathbb{E}[X]^2 && \text{(By lemmas 2 and 3)} \\
 &= \text{Var}[X] - \frac{\text{Var}[X]}{N} \\
 &= \frac{N-1}{N} \text{Var}[X]
 \end{aligned}$$

□

**Lemma 5.**  $\mathbb{E}[(x_k - \bar{x})(x_{k'} - \bar{x})] = -\frac{\text{Var}[X]}{N}$

*Proof.*

$$\begin{aligned}
 \mathbb{E}[(x_k - \bar{x})(x_{k'} - \bar{x})] &= \mathbb{E}[x_k x_{k'} + \bar{x}^2 - x_k \bar{x} - x_{k'} \bar{x}] \\
 &= \mathbb{E}[x_k] \mathbb{E}[x_{k'}] + \mathbb{E}[\bar{x}^2] - \mathbb{E}[x_k \bar{x}] - \mathbb{E}[x_{k'} \bar{x}] && \text{(By linearity and independence)} \\
 &= \mathbb{E}[X]^2 + \frac{\text{Var}[X]}{N} + \mathbb{E}[X]^2 - 2 \left( \frac{\text{Var}[X]}{N} + \mathbb{E}[X]^2 \right) && \text{(By lemmas 2 and 3)} \\
 &= \mathbb{E}[X]^2 - \frac{\text{Var}[X]}{N} + \mathbb{E}[X]^2 \\
 &= -\frac{\text{Var}[X]}{N}
 \end{aligned}$$

□

Note that the following lemmas 6-9 are generalizations of lemmas 2-5.

**Lemma 6.**  $\mathbb{E}[\hat{x}_\omega^2] = \text{Var}[X] \Omega_2 + \mathbb{E}[X]^2$

*Proof.*

$$\begin{aligned}
 \text{Var}[\hat{x}_\omega] &= \text{Var}\left[\sum_i \omega_i x_i\right] \\
 &= \sum_i \omega_i^2 \text{Var}[x_i] \\
 &= \sum_i \omega_i^2 \text{Var}[X] \\
 &= \text{Var}[X] \sum_i \omega_i^2 \\
 &= \text{Var}[X] \Omega_2
 \end{aligned}$$

And

$$\begin{aligned}
 \text{Var}[\hat{x}_\omega] &= \mathbb{E}[\hat{x}_\omega^2] - \mathbb{E}[\hat{x}_\omega]^2 \\
 &= \mathbb{E}[\hat{x}_\omega^2] - \mathbb{E}[X]^2
 \end{aligned}
 \tag{By lemma 1}$$

Which yields the desired results.  $\square$

**Lemma 7.**  $\mathbb{E}[x_i \hat{x}_\omega] = \omega_i \text{Var}[X] + \mathbb{E}[X]^2$

*Proof.*

$$\begin{aligned}
 \mathbb{E}[x_i \hat{x}_\omega] &= \mathbb{E}\left[x_i \sum_j \omega_j x_j\right] \\
 &= \omega_i \mathbb{E}[x_i^2] + \sum_{j \neq i} \omega_j \mathbb{E}[x_i] \mathbb{E}[x_j] \\
 &= \omega_i \mathbb{E}[X^2] + \sum_{j \neq i} \omega_j \mathbb{E}[X]^2 \\
 &= \omega_i \mathbb{E}[X^2] + \mathbb{E}[X]^2 (\Omega_1 - \omega_i) \\
 &= \omega_i \mathbb{E}[X^2] + \mathbb{E}[X]^2 - \mathbb{E}[X]^2 \omega_i \\
 &= \omega_i \text{Var}[X] + \mathbb{E}[X]^2
 \end{aligned}$$

$\square$

**Lemma 8.**  $\mathbb{E}[(x_i - \hat{x}_\omega)^2] = \text{Var}[X](\Omega_2 + 1 - 2\omega_i)$

*Proof.*

$$\begin{aligned}
 \mathbb{E}[(x_i - \hat{x}_\omega)^2] &= \mathbb{E}[x_i^2] + \mathbb{E}[\hat{x}_\omega^2] - 2\mathbb{E}[x_i \hat{x}_\omega] \\
 &= \mathbb{E}[X^2] + \text{Var}[X] \Omega_2 + \mathbb{E}[X]^2 - 2(\omega_i \text{Var}[X] + \mathbb{E}[X]^2) \tag{By lemmas 6 and 7} \\
 &= \mathbb{E}[X^2] + \text{Var}[X] \Omega_2 - \mathbb{E}[X]^2 - 2\omega_i \text{Var}[X] \\
 &= \text{Var}[X](\Omega_2 + 1 - 2\omega_i)
 \end{aligned}$$

$\square$

**Lemma 9.**  $\mathbb{E}[(x_i - \hat{x}_\omega)(x_j - \hat{x}_\omega)] = \text{Var}[X](\Omega_2 - \omega_i - \omega_j)$

*Proof.*

$$\begin{aligned}
 \mathbb{E}[(x_i - \hat{x}_\omega)(x_j - \hat{x}_\omega)] &= \mathbb{E}[x_i x_j + \hat{x}_\omega^2 - x_i \hat{x}_\omega - x_j \hat{x}_\omega] \\
 &= \mathbb{E}[x_i] \mathbb{E}[x_j] + \mathbb{E}[\hat{x}_\omega^2] - \mathbb{E}[x_i \hat{x}_\omega] - \mathbb{E}[x_j \hat{x}_\omega] \quad (\text{By linearity and independence}) \\
 &= \mathbb{E}[X]^2 + \text{Var}[X]\Omega_2 + \mathbb{E}[X]^2 - (\omega_i \text{Var}[X] + \mathbb{E}[X]^2) - (\omega_j \text{Var}[X] + \mathbb{E}[X]^2) \\
 &\quad \quad \quad (\text{By lemmas 6 and 7}) \\
 &= \text{Var}[X]\Omega_2 - \omega_i \text{Var}[X] - \omega_j \text{Var}[X] \\
 &= \text{Var}[X](\Omega_2 - \omega_i - \omega_j)
 \end{aligned}$$

□

**Lemma 10.**  $\mathbb{E}[\sum_i \omega_i (x_i - \hat{x}_\omega)^2] = \text{Var}[X](1 - \Omega_2)$

*Proof.*

$$\begin{aligned}
 \mathbb{E}[\sum_i \omega_i (x_i - \hat{x}_\omega)^2] &= \sum_i \omega_i \mathbb{E}[(x_i - \hat{x}_\omega)^2] \\
 &= \sum_i \omega_i \text{Var}[X](\Omega_2 + 1 - 2\omega_i) \quad (\text{By lemma 8}) \\
 &= \text{Var}[X] \left( \sum_i (\omega_i \Omega_2 + \omega_i - 2\omega_i^2) \right) \\
 &= \text{Var}[X](\Omega_2 \Omega_1 + \Omega_1 - 2\Omega_2) \\
 &= \text{Var}[X](1 - \Omega_2)
 \end{aligned}$$

□

## Proofs

*Proof.* Proof of theorem 1

Suppose that all  $w_{ij}$  are independent and prefixed weights, and assume all  $y_{ij}$  are  $\mathcal{N}(0, 1)$ .

Then we have:

$$\begin{aligned}
 S_{BB} &= \sum_{i < j} \frac{w_{ij}}{\sqrt{\sum_{i < j} w_{ij}^2}} y_{ij} \\
 &\sim \sum_{i < j} \frac{w_{ij}}{\sqrt{\sum_{i < j} w_{ij}^2}} \mathcal{N}(0, 1) \\
 &\sim \sum_{i < j} \mathcal{N}(0, \frac{w_{ij}^2}{\sum_{i < j} w_{ij}^2}) \\
 &\sim \mathcal{N}(0, \frac{\sum_{i < j} w_{ij}^2}{\sum_{i < j} w_{ij}^2}) \\
 &\sim \mathcal{N}(0, 1)
 \end{aligned}$$

□

*Proof.* Proof of theorem 2 We have  $r_\omega = \frac{\sum_{k \in I} w_k (x_k - \hat{x}_\omega)(y_k - \hat{y}_\omega)}{\sqrt{\sum_{k \in I} w_k (x_k - \hat{x}_\omega)^2} \sqrt{\sum_{k \in I} w_k (y_k - \hat{y}_\omega)^2}}$

Now note the following:

$$\begin{aligned}
 \mathbb{E} \left[ \sum_{k \in I} w_k (x_k - \hat{x}_\omega)(y_k - \hat{y}_\omega) \right] &= \sum_{k \in I} w_k \mathbb{E}[x_k - \hat{x}_\omega] \mathbb{E}[y_k - \hat{y}_\omega] \\
 &\quad \text{(By independence and linearity of expectation)} \\
 &= \sum_{k \in I} w_k (\mathbb{E}[x_k] - \mathbb{E}[\hat{x}_\omega])(\mathbb{E}[y_k] - \mathbb{E}[\hat{y}_\omega]) \\
 &= \sum_{k \in I} w_k (\mathbb{E}[X] - \mathbb{E}[X])(\mathbb{E}[Y] - \mathbb{E}[Y]) \quad \text{(By lemma 1)} \\
 &= 0
 \end{aligned}$$

Assuming we can somehow use  $\mathbb{E} \left[ \frac{A}{B} \right] = \frac{\mathbb{E}[A]}{\mathbb{E}[B]}$  that concludes the proof.

□



*Proof.* Proof of theorem 3

We have  $r^2 = \frac{Num}{Denom}$  where:

$$Num = \left[ \sum_{k \in I} (x_k - \bar{x})(y_k - \bar{y}) \right]^2$$

$$Denom = \left[ \sum_{k \in I} (x_k - \bar{x})^2 \sum_{k \in I} (y_k - \bar{y})^2 \right]$$

We have:

$$\begin{aligned} \mathbb{E}[Denom] &= \mathbb{E} \left[ \sum_{k \in I} (x_k - \bar{x})^2 \sum_{k \in I} (y_k - \bar{y})^2 \right] \\ &= \sum_{k \in I} \mathbb{E}[(x_k - \bar{x})^2] \sum_{k \in I} \mathbb{E}[(y_k - \bar{y})^2] && \text{(By independence and linearity)} \\ &= \sum_{k \in I} \frac{N-1}{N} \text{Var}[X] \sum_{k \in I} \frac{N-1}{N} \text{Var}[Y] && \text{(By lemma 4)} \\ &= (N-1)^2 (\text{Var}[X] + \text{Var}[Y]) \end{aligned}$$

$$\begin{aligned} \mathbb{E}[Num] &= \mathbb{E} \left[ \left( \sum_{k \in I} (x_k - \bar{x})(y_k - \bar{y}) \right)^2 \right] \\ &= \mathbb{E} \left[ \sum_{k \in I} (x_k - \bar{x})^2 (y_k - \bar{y})^2 + 2 \sum_{k, k' \in I, k \neq k'} (x_k - \bar{x})(x_{k'} - \bar{x})(y_k - \bar{y})(y_{k'} - \bar{y}) \right] \\ &= \sum_{k \in I} \mathbb{E}[(x_k - \bar{x})^2] \mathbb{E}[(y_k - \bar{y})^2] + 2 \sum_{k, k' \in I, k \neq k'} \mathbb{E}[(x_k - \bar{x})(x_{k'} - \bar{x})] \mathbb{E}[(y_k - \bar{y})(y_{k'} - \bar{y})] \\ &&& \text{(By linearity and independence)} \\ &= \sum_{k \in I} \frac{N-1}{N} \text{Var}[X] \frac{N-1}{N} \text{Var}[Y] + 2 \sum_{k, k' \in I, k \neq k'} (-1) \frac{\text{Var}[X]}{N} (-1) \frac{\text{Var}[Y]}{N} \\ &&& \text{(By lemmas 4 and 5)} \\ &= \left( \frac{N-1}{N} \right)^2 \text{Var}[X] \text{Var}[Y] \left( \sum_{k \in I} 1 \right) + 2 \frac{1}{N^2} \text{Var}[X] \text{Var}[Y] \left( \sum_{k, k' \in I, k \neq k'} 1 \right) \\ &= \left( \frac{N-1}{N} \right)^2 \text{Var}[X] \text{Var}[Y] (N) + 2 \frac{1}{N^2} \text{Var}[X] \text{Var}[Y] \frac{N(N-1)}{2} \\ &= \frac{(N-1)^2}{N} \text{Var}[X] \text{Var}[Y] + \frac{N-1}{N} \text{Var}[X] \text{Var}[Y] \\ &= \frac{(N-1)^2 + N-1}{N} \text{Var}[X] \text{Var}[Y] \\ &= \frac{N(N-1)}{N} \text{Var}[X] \text{Var}[Y] \\ &= (N-1) \text{Var}[X] \text{Var}[Y] \end{aligned}$$

Assuming we can somehow use  $\mathbb{E} \left[ \frac{A}{B} \right] = \frac{\mathbb{E}[A]}{\mathbb{E}[B]}$  that concludes the proof.  $\square$

*Proof.* Partial proof of conjecture 1.

We have  $r_\omega^2 = \frac{Num}{Denom}$  where:

$$Num = \left[ \sum_{k \in I} \omega_k (x_k - \hat{x}_\omega) (y_k - \hat{x}_\omega) \right]^2$$

$$Denom = \left[ \sum_{k \in I} \omega_k (x_k - \hat{x}_\omega)^2 \sum_{k \in I} \omega_k (y_k - \hat{x}_\omega)^2 \right]$$

We have:

$$\begin{aligned} \mathbb{E}[Denom] &= \mathbb{E} \left[ \sum_{k \in I} \omega_k (x_k - \hat{x}_\omega)^2 \sum_{k \in I} \omega_k (y_k - \hat{x}_\omega)^2 \right] \\ &= \mathbb{E} \left[ \sum_{k \in I} \omega_k (x_k - \hat{x}_\omega)^2 \right] \left[ \sum_{k \in I} \omega_k (y_k - \hat{x}_\omega)^2 \right] && \text{(By independence)} \\ &= \text{Var}[X](1 - \Omega_2) \text{Var}[Y](1 - \Omega_2) && \text{(By lemma 10)} \\ &= (1 - \Omega_2)^2 \text{Var}[X] \text{Var}[Y] \end{aligned}$$

And after long number crunching that we omit but will add later if desired (same idea as proof of theorem 1, using lemmas 6-9) we get:

$$\mathbb{E}[Num] = \Omega_2^3 + 2\Omega_2^2 + \Omega_2 + 4\Omega_4 - 4\Omega_2\Omega_3 - 4\Omega_3 + 2\Omega_2^2\Omega_{1,1} + 4\Omega_{1,3} + 4\Omega_{2,2} - 8\Omega_2\Omega_{1,2}$$

Where  $\Omega_{I,J} = \sum_{i \neq j} w_i^I w_j^J$

$\square$

*Proof.* Problems for proof of conjecture 2. Neither  $\mathbb{E}[(r_\alpha - r_\beta)^2] \approx \sum_i (\alpha_i - \beta_i)^2$  nor  $\mathbb{E}[r_\alpha r_\beta] \approx \sum_i \alpha_i \beta_i$  seem easy to prove.

The first makes a minus appear in the calculation which makes things very hard. Combining to common denominators makes the numerator completely nuts.

The second causes the denominator to have square roots, blocking us from continuing any further (and the numerator becomes decently more complex than it already was for conjecture 1...).

$\square$

## Results

In both subsections we used 100 runs, 100000 individuals, 100 snps per phenotype, and ran with  $h \in \{1, .5, .3, .2, .15, .1, .05\}$ .

### Previous results - before the new variance fix

Total Heritability	Regular Buhmbox Mean and Std	Continuous Buhmbox Mean and Std
1.00	-0.08, 1.00	0.23, 0.75
0.50	0.11, 1.04	0.24, 0.78
0.30	-0.08, 0.96	0.13, 0.67
0.20	0.08, 1.06	0.21, 0.74
0.15	0.12, 0.93	0.22, 0.74
0.10	-0.24, 0.96	-0.04, 0.76
0.05	0.09, 1.03	0.06, 0.69

Table 1: Independent Population

Total Heritability	Regular Buhmbox Mean and Std	Continuous Buhmbox Mean and Std
1.00	-17.46, 0.62	-17.61, 0.85
0.50	-8.36, 0.76	-8.55, 0.72
0.30	-5.18, 0.96	-5.17, 0.79
0.20	-3.11, 0.96	-3.26, 0.70
0.15	-2.57, 0.90	-2.52, 0.69
0.10	-1.68, 0.98	-1.53, 0.77
0.05	-0.66, 1.02	-0.62, 0.73

Table 2: Pleiotropic Population

Total Heritability	Regular Buhmbox Mean and Std	Continuous Buhmbox Mean and Std
1.00	39.92, 1.36	34.32, 1.13
0.50	19.41, 1.33	17.03, 0.89
0.30	11.29, 1.07	10.25, 0.79
0.20	7.55, 1.20	6.82, 0.82
0.15	5.24, 1.08	4.94, 0.78
0.10	3.73, 1.15	3.51, 0.81
0.05	1.47, 0.89	1.66, 0.69

Table 3: Heterogeneous Population

## New results - after the new variance fix

Total Heritability	Regular Buhmbox Mean and Std	Continuous Buhmbox Mean and Std
1.00	-0.03, 0.92	0.20, 0.94
0.50	0.08, 1.05	0.09, 0.93
0.30	-0.08, 0.89	0.12, 1.03
0.20	0.12, 0.92	0.27, 0.95
0.15	0.14, 0.96	0.27, 0.97
0.10	-0.07, 1.14	0.26, 1.05
0.05	-0.05, 1.01	0.16, 0.88

Table 4: Independent Population

Total Heritability	Regular Buhmbox Mean and Std	Continuous Buhmbox Mean and Std
1.00	-17.37, 0.64	-24.68, 0.99
0.50	-8.46, 0.70	-12.08, 1.01
0.30	-5.02, 0.91	-7.01, 1.00
0.20	-3.32, 0.85	-4.65, 1.02
0.15	-2.44, 0.98	-3.38, 1.08
0.10	-1.72, 1.01	-2.27, 0.98
0.05	-0.79, 0.93	-1.05, 0.85

Table 5: Pleiotropic Population

Total Heritability	Regular Buhmbox Mean and Std	Continuous Buhmbox Mean and Std
1.00	40.28, 1.34	48.94, 1.46
0.50	19.32, 1.21	24.11, 1.27
0.30	11.41, 1.22	14.45, 1.23
0.20	7.31, 1.11	9.61, 1.06
0.15	5.43, 0.99	7.26, 1.03
0.10	3.36, 0.98	4.69, 1.06
0.05	1.62, 1.01	2.44, 0.99

Table 6: Heterogeneous Population

## Remaining Problems and Challenges

I obviously still need proofs for conjectures 1 and 2. However a few other problems also present themselves.

- (1) To prove any of the theorems (or even start proving them or do anything), I had to use the following equality:

$$\mathbb{E} \left[ \frac{A}{B} \right] = \frac{\mathbb{E}[A]}{\mathbb{E}[B]}$$

Where  $A$  is more or less the squared covariance and  $B$  is the product of the two variances (since  $r = \frac{\text{Cov}[X,Y]}{\sqrt{\text{Var } X \text{Var } Y}}$  and so  $r^2 = \frac{\text{Cov}[X,Y]^2}{\text{Var } X \text{Var } Y}$ ).

Professor Pe'er suggested doing this at some point I think, and it seems to give correct results. However, it is definitely not true in general. Even if  $A$  and  $B$  are completely independent we have  $\mathbb{E} \left[ \frac{A}{B} \right] = \mathbb{E}[A] \mathbb{E}[\frac{1}{B}]$ , but I can see no reason why we would have  $\mathbb{E}[\frac{1}{B}] = \frac{1}{\mathbb{E}[B]}$  (since this is not true for almost any RV).

- (2) Even though I do not yet have all the proofs, I implemented continuous BB based on the conjecture 2 result. This correctly gave me a variance of 1 for the BB score under independence. However, the expected value of the BB score was clearly not 0. I looked at the values of various  $y_{ij}$ 's and these were all  $\mathcal{N}(0, 1)$  as expected. So the error could only come from some sort of positive covariance/correlation between the weights  $w_{ij}$  and the values of  $y_{ij}$  (rendering the result of theorem 1 inapplicable). This seems strongly problematic. I have not spent much time thinking about the problem but am unsure how to solve it.