

**Correcting For Bias In Correlation Coefficients  
Due To Intraindividual Variability**

Jeffrey Joe Nanda

Department of Mathematics  
Pomona College

Thesis submitted for the B.A Degree In Mathematics at  
Pomona College

... Spring 2007 ...

## Abstract

Correlations between variables are important in many contexts. However, when one or both of the variables exhibit intraindividual variability, conventional estimates of the Spearman rank correlation coefficient are biased towards zero. This bias towards zero is referred to as attenuation. The bias can be reduced if the mean of several measurements is used as the score rather than one measurement. Using means, however, does not eliminate the bias completely. The bias has an inverse relationship with the repeatability of the variables. In this thesis, I present an estimator for the correlation coefficient that mitigates the attenuation. The estimator is a product of the Spearman correlation coefficient and a correction factor. The correction factor is a function of within- and between-individual components of variance for each of the two traits being correlated. Simulations show that optimal sampling effort usually involves a small number of trials per individual and a large number of individuals.

*"A knowledge of statistics is like a knowledge of a foreign language or of algebra; it may prove of use at any time under any circumstances."*

A.L Bowley

# Contents

|  |            |
|--|------------|
| <b>Acknowledgments</b>   | <b>iii</b> |
| <b>1 Introduction</b>  | <b>1</b>   |
| 1.1 Introduction . . . . .   | 1          |
| <b>2 Correction for attenuation due to intraindividual variability</b> | <b>3</b>   |
| 2.1 Alternative Approaches . . . . .                                   | 3          |
| 2.2 Deriving Correction Factor . . . . .                               | 4          |
| 2.3 Simulations . . . . .  | 7          |
| 2.4 Results . . . . .  | 8          |
| 2.5 Sampling Effort . . . . .  | 11         |
| <b>3 Conclusions</b>   | <b>13</b>  |
| <b>Bibliography</b>  | <b>14</b>  |

# List of Figures

- 2.1 Distribution of  $\hat{r}$  . . . . . 9
- 2.2 Biasedness and Variation of  $\hat{r}$  with  $\rho = 0.4$  . . . . . 10
- 2.3 Biasedness and Variation of  $\hat{r}$  with  $\rho = 0.8$  . . . . . 10
- 2.4 Trade-off Between  $N_{subjects}$  and  $n_{trials}$  with  $\rho = 0.4$  . . . . . 12
- 2.5 Trade-off Between  $N_{subjects}$  and  $n_{trials}$  with  $\rho = 0.4$  . . . . . 12

# Acknowledgments

I would like to thank Pomona College Professor Johanna Hardin for her earlier work with correlation coefficients which provided interest in the subject. I would also like to thank her for being my primary contact whenever I could not understand anything and for making sure the entire project made sense. I would also like to thank May Thet Zin for proof reading the paper and providing a non-mathematical point of view.

# Chapter 1

## Introduction

### 1.1 Introduction

Correlations between various phenomena are probably one of the most widely used and frequently misused statistical tools in the natural and behavioral sciences [3]. Correlations between traits, for example, are important in helping determine functional relationships between biochemical, morphological, and whole organism traits [2]. Quantitative geneticists have always appreciated the role of phenotypic and genetic traits in determining how selection acts and how traits respond to selection. The strength of correlation coefficients between various phenomena can lead to additional research into the nature of the relationship between the variables.

Given the practical importance of correlations, accurate estimates of the strength of the linear relationships between variables are desirable. Unfortunately, conventional estimates of the Spearman rank correlation coefficient are biased towards zero whenever one or both of the variables involved in calculating a correlation coefficient exhibits intraindividual variability (i.e., the repeatability is less than one). This phenomenon is known as attenuation. Since almost all behavioral and physiological traits exhibit some degree of intraindividual variability, it is important to find an unbiased correlation coefficient. Therefore, it is safe to assume that most of the correlation coefficients reported in the literature are underestimates, on average, of the true correlation coefficients. The bias is reduced but not eliminated if the mean of several measurements is used as each individual's score rather than one measurement.

Research on finding a correction factor to mitigate the bias in the correlation coefficient has already been done for the Pearson product moment correlation coefficient. My research will focus on finding the correction factor for the Spearman rank correlation coefficient. The Spearman correlation coefficient uses ranks of measurements instead of actual measurements in calculating the correlation coefficient. The Spearman correlation coefficient is widely used in physiological and behavioral studies [5]. The chief advantage of using the Spearman rank correlation coefficient over the Pearson correlation coefficient is that the Spearman correlation coefficient reduces the effect of extreme values in calculating the correlation coefficient. In the Pearson correlation coefficient, outlying cases will have a large effect on the mean and therefore on the correlation. In a Spearman correlation, however, the extreme value is assigned a rank which is of same order of magnitude as the rest the values. Rank based correlation coefficients also eliminate disparities caused by differences in distributions of the two characteristics being correlated [9]. Using ranks, a vari-

able that follows a normal distribution can still be compared on even terms with another variable whose distribution is, for example, quite skewed. In addition, rank based correlation coefficients are also very useful if the variables being correlated have scores rather than values. A good example is when one is investigating the relationship between scores assigned to contestants in a beauty contest and contestants' heights. Since there are no units associated with beauty levels, scores are used instead. Also, sometimes, the only data available to a researcher are ranked and so is it very useful to have a rank based correlation coefficient.

In this thesis, I describe the bias that results from intraindividual variability and present an estimator that gives theoretically unbiased correlation coefficients. In section 2.1 I present alternative methods available to mitigating attenuation bias. In section 2.2 I show the theoretical derivation of the correction factor. In section 2.3 I describe the simulations I did to determine whether the correction factor mitigates the attenuation. Section 2.4 shows that the magnitude of the bias varies with the repeatability of the trait as well as the number of measurements per individual. In section 2.4, I also use simulations to evaluate the performance of my unbiased estimator versus the conventional estimator. Section 2.5 explores the allocation of effort between number of individuals versus number of measurements per individuals.

## Chapter 2

# Correction for attenuation due to intraindividual variability

### 2.1 Alternative Approaches

Conventional correlation coefficients are attenuated whenever one or both of the variables being correlated exhibit intraindividual variation. The intraindividual variation is statistically identical to measurement error. Variables that exhibit intraindividual variation have repeatabilities less than one. Repeatability is loosely defined as a statistic that describes the variation in measurements obtained with the same method on identical test material under the same test conditions (same operator, same apparatus and same environment). Variables with repeatability equal to one are those that have no measurement errors, i.e., under identical conditions, a researcher gets the same value everytime the experiment is performed e.g., number of students in a class. Variables with high repeatabilities are those that have low measurement errors and variables with low repeatabilities are those that have high measurement errors.

Researchers have always known that conventional estimates of correlation coefficients are biased if there is measurement error. Spearman [9] noted that measurement error was equivalent to rolling a ball down a rugged slope. Even though a positive association is expected between how far the ball rolls and the force with which the ball was rolled, the unevenness of the ground will confound how far the ball will roll. If a correlation coefficient is calculated between how far the ball rolled and the force with which the ball was rolled, the value obtained will certainly be less than what it should be, it will be closer to zero, hence attenuation.

Several methods have been proposed to try to eliminate the bias. Spearman suggested that researchers should use more than one measurement of the variables under study. The researcher should take several measurements, calculate their mean and then find the correlation coefficient of the means. Spearman also suggested another method in which a third variable is found which is connected to the two variables being correlated. In such a case, the unbiased correlation coefficient would be a function of the correlation coefficients of the two variables with the third variable. Using a third variable, however, he noted would sometimes overcompensate for the bias and ends up overestimating the correlation coefficient between the variables.



Some researchers favor the use of confidence sets instead of point estimates. Instead of getting a point estimate for the correlation coefficient, they calculate an interval estimate. Although this has been around for a while [4], many researchers have ignored it because of misunderstandings regarding the application of the interval estimates [4]. Using interval estimates gives boundaries of those population parameters that would most likely have produced an obtained value.

Another method which is widely used is often discussed in the context of confirmatory factor analysis [6]. In confirmatory factor analysis, which is also referred to as structural equation modeling, measurement errors are explicitly modeled and error free correlation coefficients thus calculated. This approach for obtaining measurement-error-free correlation coefficients is well known in the area of structural modeling, but it is rarely discussed within other contexts. A paper by Fan [6] showed that the confirmatory factor analysis approach yields the same results as using a correction factor.

## 2.2 Deriving Correction Factor

In deriving the correction factor, we have two variables X and Y which exhibit intraindividual variation. To model the presence of measurement error, we add the error terms  $\epsilon_{ij}$  and  $\omega_{ik}$  to the equations. Hence we have the following equations;

$$x_{ij} = \mu_{xi} + \epsilon_{ij} \quad (2.1)$$

$$y_{ik} = \mu_{yi} + \omega_{ik} \quad (2.2)$$

where;

- $i$  represents the individual
- $j, l$  are the trials
- $\mu_{xi}$  is the mean for individual  $i$  at temperature X.
- $\mu_{yi}$  is the mean of the individual  $k$  at temperature Y.

We denote the correlation between  $\bar{x}$  and  $\bar{y}$  as  $\text{cor}(\bar{x}, \bar{y})$ .

The derivation of the correction factor proceeds as following:

By definition;

$$r = \text{cor}(\bar{x}, \bar{y}) = \frac{\text{cov}(\bar{x}, \bar{y})}{\sqrt{\text{var}(\bar{x})\text{var}(\bar{y})}} \quad (2.3)$$

We know,

$$\begin{aligned} \text{cov}(\bar{x}, \bar{y}) &= \text{cov}\left(\frac{1}{n_x} \sum_{j=1}^{n_x} x_{ij}, \frac{1}{n_y} \sum_{l=1}^{n_y} y_{ik}\right) \\ &= \frac{1}{n_x} \frac{1}{n_y} \sum_{j=1}^{n_x} \sum_{l=1}^{n_y} \text{cov}(x_{ij}, y_{ik}) \\ &= \frac{1}{n_x} \frac{1}{n_y} \sum_{j=1}^{n_x} \sum_{l=1}^{n_y} \text{cov}(\mu_{xi} + \epsilon_{ij}, \mu_{yi} + \omega_{ik}) \\ &= \frac{1}{n_x} \frac{1}{n_y} \sum_{j=1}^{n_x} \sum_{l=1}^{n_y} (\text{cov}(\mu_{xi}, \mu_{yi}) + (\text{cov}(\mu_{xi}, \omega_{ik}) + \text{cov}(\mu_{yi}, \epsilon_{ij}) + \text{cov}(\epsilon_{ij}, \omega_{ik}))) \end{aligned}$$

But we also know that

$$\begin{aligned} \text{cov}(\mu_{xi}, \omega_{ik}) &= \text{cov}(\mu_{yi}, \epsilon_{ij}) \\ &= \text{cov}(\epsilon_{ij}, \omega_{ik}) \\ &= 0 \text{ (by independence)} \end{aligned}$$

Canceling the zero terms we get,

$$\begin{aligned} \text{cov}(\bar{x}, \bar{y}) &= \frac{1}{n_x} \frac{1}{n_y} \sum_{j=1}^{n_x} \sum_{l=1}^{n_y} \text{cov}(\mu_{xi}, \mu_{yi}) \\ &= \text{cov}(\mu_{xi}, \mu_{yi}). \end{aligned} \tag{2.5}$$

In addition, we also know that by definition,  $\rho$  is the true correlation coefficient of the true means of the variables X and Y, i.e.,  $\mu_{xi}$  and  $\mu_{yi}$ . Hence,

$$\rho = \text{cor}(\mu_{xi}, \mu_{yi}) = \frac{\text{cov}(\mu_{xi}, \mu_{yi})}{\sqrt{\text{var}(\mu_{xi})\text{var}(\mu_{yi})}} \tag{2.6}$$

By simple algebra,

$$\text{cov}(\mu_{xi}, \mu_{yi}) = \rho \sqrt{\text{var}(\mu_{xi})\text{var}(\mu_{yi})}. \tag{2.7}$$

To find  $\text{var}(\bar{x})$ , we proceed as follows:

$$\begin{aligned} \text{var}(\bar{x}) &= \text{var}\left(\frac{1}{n_x} \sum_{j=1}^{n_x} x_{ij}\right) \\ &= \frac{1}{n_x^2} \text{var}\left(\sum_{j=1}^{n_x} \mu_{xi} + \epsilon_{ij}\right) \\ &= \frac{1}{n_x^2} \text{var}\left(n_x \mu_{xi} + \sum_{j=1}^{n_x} \epsilon_{ij}\right) \\ &= \frac{n_x^2}{n_x^2} \text{var}(\mu_{xi}) + \frac{1}{n_x^2} \text{var}\left(\sum_{j=1}^{n_x} \epsilon_{ij}\right) \text{ (by independence)} \\ &= \text{var}(\mu_{xi}) + \frac{1}{n_x} \text{var}(\epsilon_{ij}) \end{aligned}$$

Hence,

$$\text{var}(\bar{x}) = \text{var}(\mu_{xi}) + \frac{1}{n_x} \text{var}(\epsilon_{ij}) \tag{2.8}$$

Similarly,

$$\text{var}(\bar{y}) = \text{var}(\mu_{y_k}) + \frac{1}{n_y} \text{var}(\omega_{ik}) \quad (2.9)$$

Substituting (2.7), (2.8), and (2.9) into (2.3) we get

$$\begin{aligned} r &= \text{cor}(\bar{x}, \bar{y}) \\ &= \rho \frac{\sqrt{\text{var}(\mu_{xi})\text{var}(\mu_{yi})}}{\sqrt{(\text{var}(\mu_{xi}) + \frac{1}{n_x} \text{var}(\epsilon_{ij}))(\text{var}(\mu_{yi}) + \frac{1}{n_y} \text{var}(\omega_{ik}))}} \\ &= \rho \sqrt{\frac{\sigma_{a,X}^2 \sigma_{a,Y}^2}{(\sigma_{a,X}^2 + \frac{\sigma_{w,X}^2}{n_x})(\sigma_{a,Y}^2 + \frac{\sigma_{w,Y}^2}{n_y})}} \\ &= \rho \sqrt{\frac{\sigma_{a,X}^2}{(\sigma_{a,X}^2 + \frac{\sigma_{w,X}^2}{n_x})} \frac{\sigma_{a,Y}^2}{(\sigma_{a,Y}^2 + \frac{\sigma_{w,Y}^2}{n_y})}} \end{aligned}$$

Which gives

$$\text{cor}(\bar{x}, \bar{y}) = \rho \sqrt{\frac{\sigma_{a,X}^2}{(\sigma_{a,X}^2 + \frac{\sigma_{w,X}^2}{n_x})} \frac{\sigma_{a,Y}^2}{(\sigma_{a,Y}^2 + \frac{\sigma_{w,Y}^2}{n_y})}} \quad (2.10)$$

Where,

- $\sigma_{a,X}^2$  is the true variation among the individuals at temperature X.
- $\sigma_{a,Y}^2$  is the true variation among the individuals at temperature Y.
- $\sigma_{w,X}^2$  is the true variation in the values within an individual at temperature X.
- $\sigma_{w,Y}^2$  is the true variation in the values within an individual at temperature Y.
- $n_x$  is the number of measurements (trials) per individual at temperature X.
- $n_y$  is the number of measurements (trials) per individual at temperature Y.

Before we proceed to show the results of the theoretical derivation, it is necessary at this point to define the 4 different correlation coefficients in this thesis.

- $\rho = \text{cor}(\mu_{xi}, \mu_{yi})$ , here there are an infinite number of subjects and an infinite number of trials per individual.
- $r = \text{cor}(\mu_{xi}, \mu_{yi})$ , here we have an infinite number of subjects, but a finite number of trials.
- $\hat{r} = \text{cor}(\bar{x}, \bar{y})$ , this is the conventional correlation coefficient with finite number of subjects and a finite number of trials per individual.
- $r_{corrected} = (\hat{r})(\text{correction factor})$ .

To find the correction factor, we notice that  $r = \rho z$  where z is the term under the square root in equation (2.10).

Hence we now have the relationship  $\rho = \frac{z}{a}$ . Based on the equations above and a little algebraic manipulation of  $z$ , we get,

$$\rho = r \sqrt{\left(1 + \frac{\sigma_{w,X}^2}{n_x \sigma_{a,X}^2}\right) \left(1 + \frac{\sigma_{w,Y}^2}{n_y \sigma_{a,Y}^2}\right)} \quad (2.11)$$

If we are using values from samples, we get the following equation,

$$r_{corrected} = \hat{r} \sqrt{\left(1 + \frac{s_{w,X}^2}{n_x s_{a,X}^2}\right) \left(1 + \frac{s_{w,Y}^2}{n_y s_{a,Y}^2}\right)} \quad (2.12)$$

Where,

- $s_{a,X}^2$  is the sample variation among the individuals at temperature X.
- $s_{a,Y}^2$  is the sample variation among the individuals at temperature Y.
- $s_{w,X}^2$  is the sample variation in the values within an individual at temperature X.
- $s_{w,Y}^2$  is the sample variation in the values within an individual at temperature Y.

The correction factor is the term under the square root sign. This formula has been presented by several authors [2]. Equation (2.10) shows that if intraindividual variability exists, the conventional estimates of the correlation coefficient is biased toward zero because the term under the square root sign is always less than one. Ideally, the scale factor "z" is one which would make  $r = \rho$ . The magnitude of the bias decreases as the number of measurements per individual increases and also as the repeatability (which is given by  $\frac{\sigma_{w,Y}^2}{\sigma_{a,Y}^2}$ ) increases.

## 2.3 Simulations

Using equation (2.12), an unbiased estimator for the correlation coefficient is given by

$$r_{corrected} = \hat{r} \sqrt{\left(1 + \frac{s_{w,X}^2}{n_x s_{a,X}^2}\right) \left(1 + \frac{s_{w,Y}^2}{n_y s_{a,Y}^2}\right)} \quad (2.13)$$

where  $\hat{r}$  is the conventional correlation coefficient calculated using mean values ( $\bar{x}, \bar{y}$ ) for each individual at each of the temperatures X and Y obtained from the sample. Also, the population variance terms  $\sigma^2$  are replaced with sample variance terms  $s^2$ . This was done because we do not always know the true population values, but we can estimate the population values from the samples. Equation (2.13) can also be written in terms of repeatabilities [2];

$$r_{corrected} = \hat{r} \sqrt{\left(1 + \frac{1 - r_{i,X}}{n_x r_{i,X}}\right) \left(1 + \frac{1 - r_{i,Y}}{n_y r_{i,Y}}\right)} \quad (2.14)$$

Where;

- $r_{i,X}$  is the repeatability of trait X
- $r_{i,Y}$  is the repeatability of trait Y

Equation (2.14) can be used to obtain unbiased estimates of correlations from published work that provided results for repeatabilities. To determine whether the correction factor worked, I simulated data on locomotor performance of the lizard *Sceloporous graciosus* based on data from Adolph [1]. In the experiment, there were 21 lizards, and each lizard was raced around a track 20 times at 20° C and 20 times at 35° C. Hence for the experiment,  $n_x = n_y = 20$ . Similarly, in the simulation, there were 21 "lizards" and each lizard was "raced" around a track 20 times at each of the temperatures and all the times were then ranked. A Spearman correlation coefficient was calculated to determine the relationship between the ranks of the times at each of the temperatures. The samples were drawn from a bivariate normal distribution with known correlation coefficients ( $\rho = 0.4$  and  $0.8$ ). I examined the effects of the number of lizards ( $N_{subjects} = 2, 5, \text{ and } 10$ ) and repeatability ( $r_i = 0.2, 0.5 \text{ and } 0.8$ ) on the distribution of the corrected Spearman correlation coefficients. The number of measurements per lizard at each of the temperatures X and Y was set to be equal ( $n_x = n_y = n_{trials}$ ) and varied ( $n_{trials} = 2, 5 \text{ and } 10$ ) to examine the effect of the change in the number of the measurements on the distribution of the corrected spearman correlation coefficient. As shown in equations (2.1) and (2.2), normally distributed error terms were added and adjusted to yield the desired repeatability. For each combination of parameter values, I drew 5000 independent samples and from each sample calculated four different estimates of the correlation coefficient.

(1) uncorrected for bias, but without adding measurement error terms (i.e.  $\sigma_{w,x} = 0$ ) as a check for the simulation procedure. This should on average, yield an unbiased estimate because the repeatabilities equal 1.

(2) uncorrected for bias. This is the conventional Spearman correlation coefficient. On average, this should yield a biased estimate.

(3) corrected for bias using equation (2.10) using known population variance components. This should yield an unbiased estimate on average.

(4) corrected for bias using equation (2.10) using variance components estimated from the sample. This should also yield an unbiased estimate.

Simulations were run using the statistical language R.

## 2.4 Results

Figure (2.1) shows the distribution of sample correlation coefficients for a single set of parameter values. In this example, I used 10,000 runs. Several features are worth noting in the plot. The uncorrected correlation coefficients were biased towards zero as predicted by the theory, whereas the corrected values were generally unbiased. Another interesting feature is that the distributions differ in variance. The variance of the unbiased estimates (labeled w/o error) is entirely due to sampling of individuals because these runs have zero intraindividual variability (repeatability= 1). The other three distributions exhibit variance due to both sampling of individuals and sampling values within individuals (intraindividual variability). The greater sampling variance of the unbiased estimators compared to the conventional estimator is due to the fact that the unbiased estimate involves multiplying the biased estimate by a multiplicative factor. It is also interesting to note that the variance of the corrected correlation coefficients are about the same indicating that estimating the variance parameters does not increase the sampling variance extensively. Also, the graph showing the distribution of the corrected correlation coefficients (labeled ANOVA adj) has a bump at the end of the distribution. The bump is due to simulation coding. The correlation coefficients were forced to be less than one which

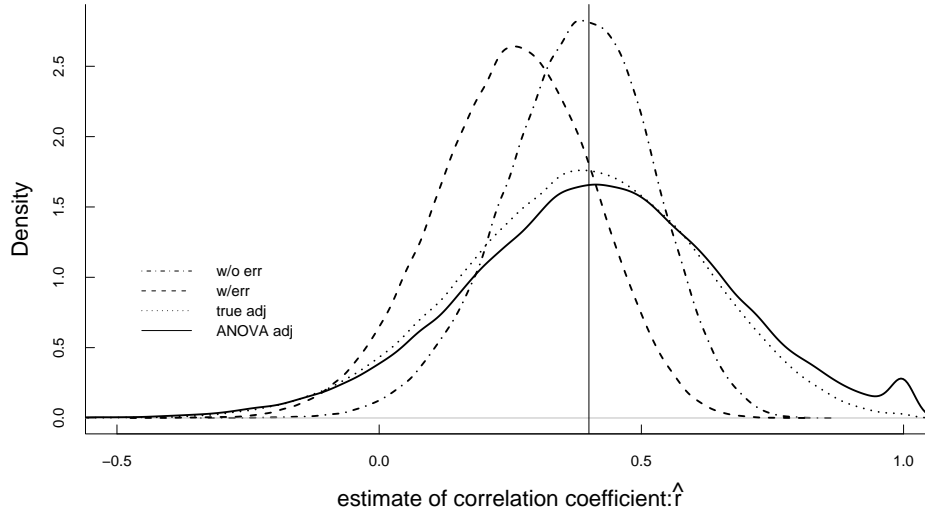


Figure 2.1: Distribution of  $\hat{r}$

ended up squashing the distribution causing the bump.

Figures 2.2 and 2.3 summarize the simulation results for all the combinations of the parameters used. The difference between the two graphs is that figure 2.2 has  $\rho = 0.4$  and figure 2.3 has  $\rho = 0.8$ . In figures 2.2 and 2.3 each boxplot represents 90% of the distribution of the corrected correlation coefficients. The bottom line is the fifth percentile, the bold middle line represents the median and the top line is the 95<sup>th</sup> percentile. The first boxplot represents 90% of the distribution of the conventional correlation coefficients when  $n_{trials} = 2$ ,  $\rho = 0.4$ ,  $N_{subjects} = 20$  and the repeatabilities ( $r_{i,X}$  and  $r_{i,Y} = 0.2$ ). The two boxplots next to the first one represent the distributions of corrected correlation coefficients using sample variance terms and population variance terms respectively but with the same values for  $n_{trials}$ ,  $\rho$ ,  $r_{i,X}$  and  $r_{i,Y}$ . The three box plots next to the first three represent the distribution of the correlation coefficients mentioned above but with  $r_{i,X}$  and  $r_{i,Y} = 0.5$

All the simulations confirmed that the conventional correlation coefficient is biased and that the correction factor eliminated the bias. The simulations also confirmed various other interesting facts. For example, the magnitude of the bias was greatest for the smallest number of samples per individual and for the lowest repeatabilities. In some of the cases, especially where the true correlation coefficient  $\rho$  was 0.8, the entire boxplot fell below the true value. When  $\rho = 0.4$ , there

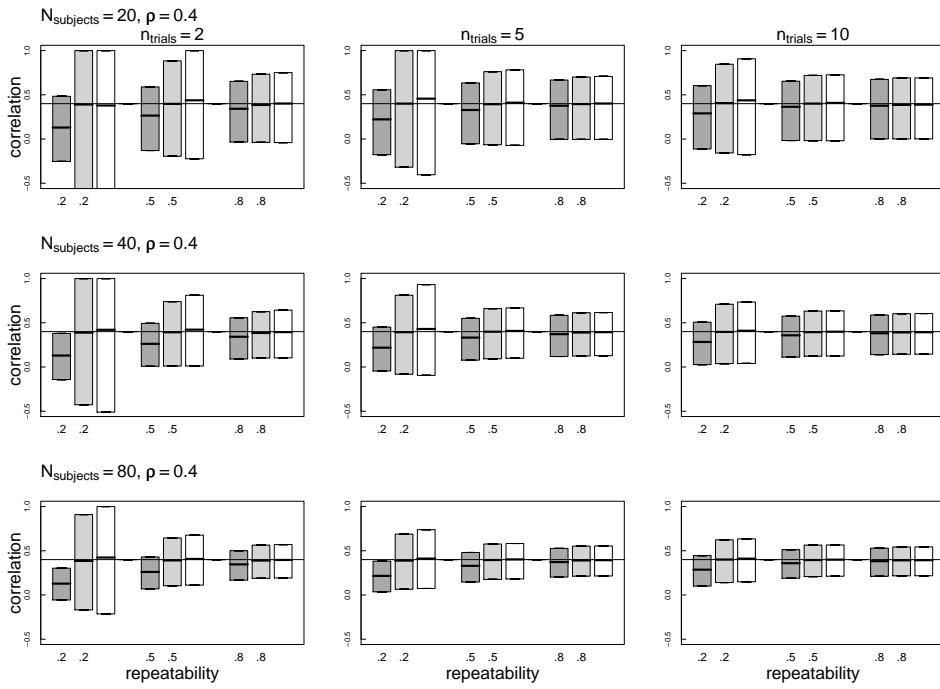


Figure 2.2: Biasedness and Variation of  $\hat{r}$  with  $\rho = 0.4$

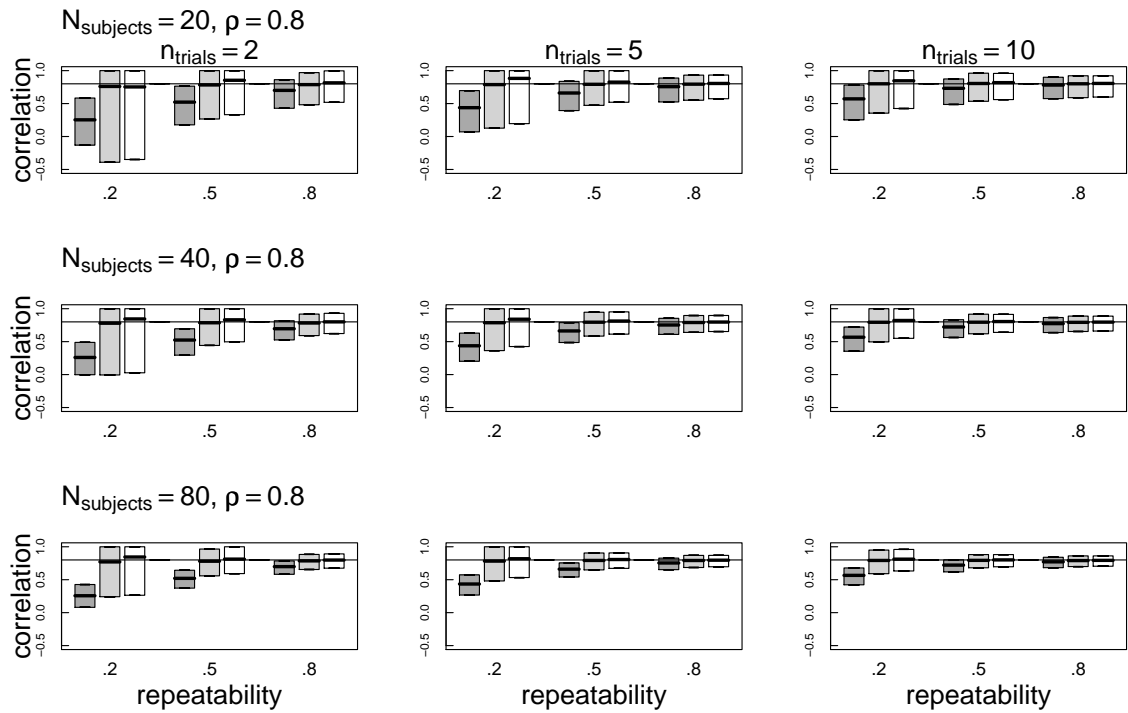


Figure 2.3: Biasedness and Variation of  $\hat{r}$  with  $\rho = 0.8$

were fewer boxplots that fell entirely below the true value even though most of the distribution of the correlation coefficients lied below the true value for lower sample sizes and lower repeatabilities. Interestingly, the attenuation is not reduced by increasing the number of subjects. Increasing number of subjects, however, reduces the variance of the distribution of the correlation coefficients.

## 2.5 Sampling Effort

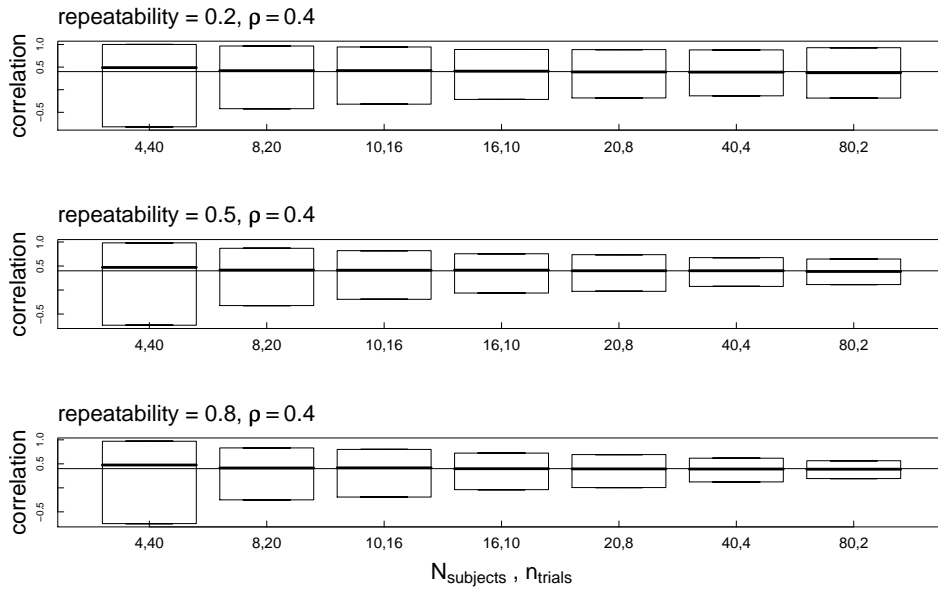
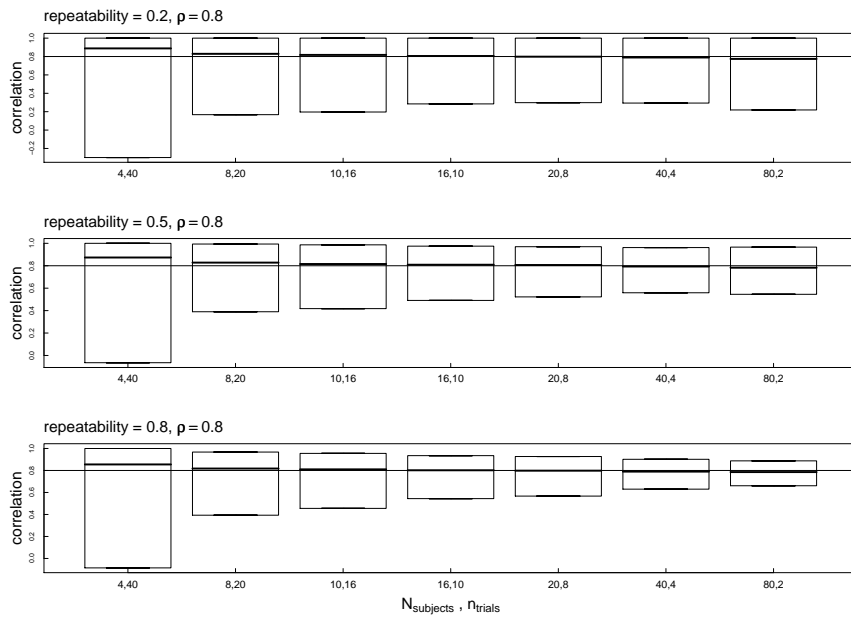
The unbiased estimator for the correlation coefficient requires at least two measurements per individual so we can estimate the within-and among-variance components used in equation (3). As noted before, having more than two measurements per individual decreases the bias and should also give more accurate estimates of the mean as well as the variance components. However, increasing the number of individuals reduces the variance of the distribution of the corrected correlation coefficients, which gives more accurate estimates. So there is a trade-off between the number of trials and the number of individuals.

I ran simulations to explore the trade-off between  $N_{subjects}$  and  $n_{trials}$  and to determine how varying them would affect the distributions of the estimates of  $\rho$ . I assumed that a researcher could only take 160 measurements and had to choose between more lizards ( $N_{subjects}$ ) and more runs per lizard ( $n_{trials}$ ). The number of lizards ranged from 4 to 80 and the number of trials per lizard ranged from 2 to 40. I then obtained 5000 samples for each allocation and for each possible combination of repeatabilities and true correlations. Assuming the researcher could only take 160 measurements was done because, sometimes obtaining measurement is very difficult or dangerous. Also, sometimes, there are considerable costs involved in running the equipment used to take the measurements and one would like to take measurements in the most efficient manner.

Figures 2.4 and 2.5 show the effect of the trade-off between number of individuals versus number of trials per individual on the variance of the unbiased correlation coefficient (corrected using sample variance components and labeled ANOVA adj in figure 2.1) obtained via simulations. Figure 2.4 shows the variation in the corrected correlation coefficients when  $\rho = 0.4$  and figure 2.5 shows the variation in the corrected correlation coefficients when  $\rho = 0.8$ .

Each bar shows the median and the central 90% of the sample correlation coefficients. The total number of observations is fixed at 160. Results from the simulations show that the estimates are unbiased since the median lies on the true correlation coefficients value and the distributions are symmetric. In the (4,40) cases, the median lines lie above the true correlation coefficient lines but the distributions are not symmetric so the estimates are not biased either. For low repeatabilities in both  $\rho = 0.4$  and 0.8, it seems that the variance is lowest for some intermediate allocation value. In this simulation, the variance was lowest when we had 40 subjects and 4 trials per subject. In higher repeatability cases however, the variance monotonically decreases as the number of individuals increases. These results suggest that if we want to minimize the sampling variance of the correlation coefficients, we have to determine the repeatability of the traits under investigation. If the traits have high repeatabilities, the most precise estimates are obtained when we have the largest number of individuals and only 2 trials per individual.



Figure 2.4: Trade-off Between  $N_{\text{subjects}}$  and  $n_{\text{trials}}$  with  $\rho = 0.4$ 

## Chapter 3

# Conclusions

The Spearman Correlation coefficients between individual mean values are biased towards zero when one or both the traits exhibit intraindividual variability. Theoretically, a multiplicative factor can be found which eliminates the bias. The correction factor is a function of among- and within-variance components of the measurements. It can also be expressed as a function of repeatabilities of the traits. Thus, most uncorrected correlation coefficients reported in literature on physiological traits are biased towards zero. The bias is however, easy to correct as long as the within- and among-individual variance components are known or are able to be calculated. The corrected correlation coefficients have higher variance than the conventional correlation coefficient estimator. This variance can be reduced by increasing the number of individuals in the study, but that does not reduce the bias. Thus, depending on the precision required by the researcher and the repeatability of the traits, an optimal trade off between number of individuals and number of trials can be found. Because two measurements per individual are enough to calculate variance components needed to eliminate the bias, optimal design for sampling effort usually requires a few measurements per individual (at least two) and the largest number of individuals feasible. The correction factor is then used to eliminate the bias.

# Bibliography

- [1] Adolph, S.C. (1987). Physiological and behavioral ecology of the lizards *Sceloporous occidentalis* and *Sceloporous garciosus*. Dissertation, University Of Washington, Seattle, WA
- [2] Adolph, S.C., Johanna Hardin. Estimating Phenotypic Correlation; Correcting for Bias Due to Intra-individual Variability. *Functional Biology* 21 (2007): 178-184
- [3] Carroll, J.B. The nature of the data, or how to choose a correlation coefficient. *Psychometrika* 26(1961): 347-372
- [4] Charles, Eric, P. The correction for attenuation due to measurement error:clarifying concepts and creating confidence sets. *Psychological Methods* 10(2005): 206-226
- [5] Chen, Peter,Paula Popovich. *Parametric and nonparametric measures*, Sage Publications, Thousand Oaks, CA. 2002
- [6] Fan, Xitao. Two approaches for correcting correlation attenuation caused by measurement error: implicatons for research practice. *Educational and Psychological Mesurement* 63(2003):915-930
- [7] Kutner,Michael, Christopher J. Nachtsheim, John Neter, William Wasserman. *Applied Linear statistical Models*, McGraw Hill Irwin Publishers, 2005.
- [8] Rosner, B., and W. C. Willett. Interval Estimates for correlation Coefficients corrected for within person variation implications for study design and hypothesis testing,*American Journal of Epidemiology* 127(1988): 377-386.
- [9] Spearman,Carl. The proof and measurement of association between two things, *American Journal of Psychology* 15(1904): 72-101.