# TO POOL OR NOT TO POOL:
# A COMPARISON BETWEEN TWO
# COMMONLY USED TEST STATISTICS

Mouchumi Bhattacharyya

Department of Mathematics
University of the Pacific
3601 Pacific Avenue, Stockton, CA 95211, USA

**Abstract:**   When it comes to testing hypotheses regarding two population means, the most commonly used test is the two-sample t-test. There are two versions of this test, one is used when the variances of the two populations are equal (the pooled test) and the other one is used when the variances of the two populations are unequal (the unpooled test). The pooled test seems to have fallen into some disfavor because of its 'claimed' sensitivity to departures from the assumptions of equal population variances. Through a simulation study, this paper demonstrates that although both the pooled and the unpooled test underperform at times in their allocated settings, the overall performance of the pooled t-test is significantly superior to that of the unpooled t-test.

**Key Words:**   pooled t-test, equal variance

## 1. Introduction

The main goal of this paper is to investigate the two-sample t-test under two different scenarios - first, when the two populations have the same variance and secondly, when the two populations have different variances. When the

two populations have or assumed to have the same variance, the test statistic $t = \frac{(\overline{x_1}-\overline{x_2})-(\mu_1-\mu_2)}{s_p\sqrt{\frac{1}{n_1}+\frac{1}{n_2}}}$ is used to test for the equality of the two population means. In this test statistic, $n_1$, $\overline{x_1}$, and $n_2$, $\overline{x_2}$ are the sample sizes and the sample means of the first and the second sample respectively, whereas $\mu_1$ and $\mu_2$ are the corresponding *population* means. Also, the quantity $s_p^2$ is the pooled sample variance which is defined by $s_p^2 = \frac{(n_1-1)s_1^2+(n_2-1)s_2^2}{n_1+n_2-2}$, where $s_1^2$ and $s_2^2$ are the sample variances of the first and the second sample respectively. Given the two populations are normal, the t-statistic mentioned above indeed has an exact t-distribution with $n_1 + n_2 - 2$ degrees of freedom. Note that if the null hypothesis $H_0 : \mu_1 = \mu_2$ is indeed true, i.e., if the two populations, in fact, have the same means, the above test statistic reduces to $t = \frac{(\overline{x_1}-\overline{x_2})}{s_p\sqrt{\frac{1}{n_1}+\frac{1}{n_2}}}$. This test statistic, from now on, will be referred to as TS-pooled. Given the two populations are normal and under the assumption that the two populations have *different* variances, the test statistic used to test for equality of the two population means is $t^* = \frac{(\overline{x_1}-\overline{x_2})-(\mu_1-\mu_2)}{\sqrt{\frac{s_1^2}{n_1}+\frac{s_2^2}{n_2}}}$. Under the assumption of normal populations, this test statistic is known to have an *approximate* t-distribution with $\nu = \frac{\left(\frac{s_1^2}{n_1}+\frac{s_2^2}{n_2}\right)^2}{\frac{\left(s_1^2/n_1\right)^2}{n_1-1}+\frac{\left(s_2^2/n_2\right)^2}{n_2-1}}$ degrees of freedom. We use the notation $t^*$ to indicate that this is an approximate, and not an exact, t-distribution. (Note that if the two populations were indeed normal, the statistic $z = \frac{(\overline{x_1}-\overline{x_2})-(\mu_1-\mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1}+\frac{\sigma_2^2}{n_2}}}$ has a standard normal distribution where $\sigma_1^2$ and $\sigma_2^2$ are the two population variances). As before, if the null hypotheisis of equal population menas were indeed true, the above test statistic will reduce to $t^* = \frac{(\overline{x_1}-\overline{x_2})}{\sqrt{\frac{s_1^2}{n_1}+\frac{s_2^2}{n_2}}}$. This test statistic, from now on, will be referred to as TS-unpooled.

The most common practice to decide between TS-pooled and TS-unpooled is to perform a hypothesis test for the equality of the two population variances. A more informal way to decide between TS-pooled and TS-unpooled is to compare the observed *sample* standard deviations. If the ratio of the two standard deviations fall between $\frac{1}{2}$ and 2 (or if one standard deviation is within two times of the other standad deviation), TS-pooled is used to test for equality of the two population means. If the ratio of the two standard deviations fall outside of the interval $\left(\frac{1}{2}, 2\right)$, then TS-unpooled is used to test for equality of the two population means.

The pooled t-test or TS-pooled, which is the theoretically the correct t-

statistic, has fallen into some disfavor because of its 'claimed' sensitivity to departures from the assumtions of equal population variances (Peck, Olsen, & Devore). We use a simulation study to disprove this claim. The study consists of 240 comparisons of the two test statistics. For the sake of simplicity, we wil describe here the basics of one such comparison. We will also introduce a few terms which we will be using throughout the paper. For a single comparison of the two test statistics, we draw two independent random samples from two simulated populations. The two populations may (or may not) have equal means and/or equal variances. For a particular comparison, if the two populations indeed have equal variances, we designate TS-pooled to be the 'correct' test statistic. Similarly if the two populations have unequal variances, we designate TS-unpooled to be the 'correct' test statistic. Once the two independent samples from the two simulated populations are drawn, we perform the test of hypothesis of equality of two means using **both** TS-pooled and TS-unpooled. Since the tests are conducted on samples from known populations, we record the conclusions of both TS-pooled and TS-unpooled as correct or incorrect. Furthermore, we label one of the two test statistics as the 'better' one if the p-value corresponding to that test-statistic is closer to the correct conclusion (unless, obviously, both test statistics have exactly the same p-value). For example, if the two populations, where the two samples were drawn from, indeed had the same mean, then the test-statistic that yielded the bigger p-value is labeled as the 'better' one, whereas if the two populations, where the two samples were drawn from, had unequal means, then the test-statistic that yielded the smaller p-value is labeled as the 'better' one. In addition, we label the test statistic that yielded the p-value which is farther away (when compared to each other) from the correct conclusion to be the 'underperformer'. Also, if the two samples were drawn from two populations with the same variance, we refer to it as the "equal-variance setting", whereas if the two samples were drawn from two populations with unequal variances, we refer to it as the "unequal-variance setting". Although we reveal, in section 3, the number of times each test statistic arrives at the correct conclusion, it is important to note that as far as the comparisons are concerned, we are strictly interested in finding the 'better' one (or the 'underperformer') of the two.

## 2. Analysis

As mentioned in the last section, we compare TS-pooled and TS-unpooled in 240 simulated tests. We divide the study into the following four phases:

1. When the two populations have different means and different variances.

2. When the two populations have different means but the same variance.

3. When the two populations have the same mean but different variances.

4. When the two populations have the same mean and the same variance.

Within each of the above four phases, we make 60 comparisons of the the two test statistics. The details of the results of each phase will be given in the following section. Here we present a summary of the overall results. Of the 120 comparisons conducted under the unequal-variance setting (phases 1 and 3 above), i.e, a setting in which TS-unpooled is the 'correct' test statistic and hence is expected to perform 'better' than TS-pooled, we found TS-unpooled to perform 'better' than TS-pooled in 56 cases (47%), TS-pooled performed 'better' than TS-unpooled in 48 cases (40%), and the two test statistics performed equally well in 16 cases (13%). Similarly, of the 120 tests conducted under the equal-variance setting (phases 2 and 4 above), i.e, a setting in which TS-pooled is the 'correct' test statistic and hence is expected to perform 'better' than TS-unpooled, we found TS-pooled to perform 'better' than TS-unpooled in 60 cases (in 50% cases), TS-unpooled performed 'better' than TS-pooled in 35 cases (in 29% cases) and the two test statistics performed equally well in 25 cases (in 21% cases). The following table shows these results.

| **Unqual**-variance setting: Sample size 120 |
| --- |
| # of timesTS-pooled was 'better'= 48 (or 40%) |
| # of timesTS-unpooled was 'better'= 56 (or 47%) |
| # of ties= 16 (or 13%) |

| **Equal**-variance setting: Sample size 120 |
| --- |
| # of timesTS-pooled was 'better'= 60 (or 50%) |
| # of timesTS-unpooled was 'better'= 35 (or 29%) |
| # of ties= 25 (or 21%) |

A closer look at the above data reveals that there is no clear 'winner' between the two test statistics, TS-pooled and TS-unpooled, in the sense that they both 'underperform' in their allocated settings more than one would have hoped. Yet, from the first three conclusions below, we find the overall performance of TS-pooled to be significantly superior to the overall performance of TS-unpooled when the two test statistics are compared to each other with respect to their performances at their allocated settings. Although the first

three conclusions below appear to be roughly the same, they are presented from different perspectives.

**Conclusion 1.** *When compared to each other with respect to their allocated settings, TS-pooled was found to be more reliable than TS-unpooled. In other words, the performance of TS-pooled in the equal-variance setting is significantly superior to the performance of TS-unpooled in the unequal-variance setting (p-value = 0.0336). This is based on the fact that in the equal-variance setting, the number of times TS-pooled performed as good or better than TS-unpooled was 85 (or 71%), whereas in the unequal-variance setting, the number of times TS-unpooled performed as good or better than TS-pooled was 72 (or 60%).*

**Conclusion 2.** *One can make the same conclusion above from another perspective which is that TS-unpooled underperforms in its allocated setting (unequal variance setting) significantly more frequently than TS-pooled underperforms in its allocated setting (equal variance setting) (p-value = 0.0336). This is based on the fact that in our samples, the number of times TS-unpooled underperformed in the unequal-variance setting was 48 (or 40%), whereas the number of times TS-pooled underperformed in the equal-variance setting was only 35 (or 29%).*

**Conclusion 3.** *Yet another perspective, and probably a more interesting way of interpreting the above two conclusions is to conclude that when compared how the two test statistics perform when the underlying variance assumptions of their allocated settings were not met, TS-pooled is significantly 'better' than TS-unpooled (p-value = 0.0336). This is due to the fact that in the unequal-variance setting, TS-pooled performed better than TS-unpooled 48 times (40%), whereas in the equal-variance setting TS-unpooled performed better than TS-pooled only 35 times (29%).*

**Conclusion 4.** *When the performance of TS-unpooled was compared with itself under the two settings - equal-variance and nequal-variance, we found that the proportion of times TS-unpooled performs 'better' than TS-pooled in the it's allocated setting (unequal-variance setting) is significantly higher than the proportion of times TS-unpooled performs 'better' than TS-pooled in the equal-variance setting (p-value=0.0021). This is due to the fact that in the unequal-variance setting the number of times TS-unpooled performed 'better' than TS-pooled was 56 (47%), whereas in the equal-variance setting the number of times TS-unpooled performed 'better' than TS-pooled was 35 (29%). In a similar comparison when the performance of TS-pooled was compared with itself under the two settings - equal-variance and unequal-variance, we found*

*that the proportion of times TS-pooled performs 'better' than TS-unpooled in the equal variance setting was **not** found to be significantly higher than the proportion of times TS-pooled performs 'better' than TS-unpooled in the unequal variance setting (p-value = 0.0594). This is due to the fact that in the equal-variance setting the number of times TS-pooled performed 'better' than TS-unpooled was 60 (50%), whereas in the unequal-variance setting the number of times TS-pooled performed 'better' than TS-unpooled was 48 (40%) (p-value=0.0594). It is important to note that this conclusion (conclusion 4) would reverse between TS-unpooled and TS-pooled (with the exact same p-values) if we were to change the criterion 'better' to 'as well or better'.*

## 3. Data and Results

Here we elaborate on how the four phases of the study were conducted. For each phase of the study, two normal populations, each with ten thousand observations, were generated using the statistical software package Minitab. Although normality of the populations were not required for this study as the sample sizes were relatively large (central limit theorem would guarantee the desired results), normal populations were generated for maximum flexibility. Since the purpose of this paper was to simply comapre TS-pooled and TS-unpooled, we chose to deal with normal populations **and** large samles to minimize any type of bias. Also, in order to cover various scenarios, the population means and/or the standard deviations were varied within each phase of the study, as well as the sample sizes drawn from the two populations. Once the sizes of the two samples from the two populations were determined within a particular phase, several sets of samples were drawn repeatedly (at least fifteen sets). For each set within a phase, two-sample t-tests were carried out, using both TS-pooled and TS-unpooled, to test the hypothesis about equality of the two population means. Considering that the samples were drawn from known populations, the conclusions of the hypotheses tests could be verified as correct or incorrect. The purpose of doing the study in several phases, by changing the population means and/or the standard deviations, is to investigate how well the two test statistics perform in terms of arriving at the correct conclusions based on the proximity of the two population means from each other. In this section we present the results of each phase in the following format.

- We first lay out the two populations with their respective means and standard deviations.

- We specify the sample sizes taken from each population.

- We specify the number of *repeated* samples taken.

- Based on the standard deviations of the two *known* populations (equal or unequal), we designate one of the two test statistics (TS-pooled or TS-unpooled) as the 'correct' test statistic in the sense that the 'correct' test statistic is expected to lead to a 'better' conclusion compared to the other test statistic.

- We reveal the number of times each of the two test statistics arrive at the correct conclusion.

- We finally point out how often each of the test statistics perform better (and as well as) compared to the other one, in terms of having a p-value "closer" to the correct conclusion.

Note that the significance level was taken to be 0.05 all of our conclusions.

### Phase 1: The two populations have different means and different variances
### (Stages 1 through 4 below)

**Stage 1:**

| Population 1: $\mu_1 = 75$    $\sigma_1 = 15$ | Population 2: $\mu_2 = 70$    $\sigma_2 = 7$ |
|---|---|
| $\downarrow$ | $\downarrow$ |
| $n_1 = 50$ | $n_2 = 60$ |

- Number of repeated samples $= 15$

- **Correct** Test Statistic: **TS-unpooled**

- Number of times TS-unpooled was able to arrive at the correct conclusion (i. e., reject the null hypothesis of equal mean) $= 7$

- Number of times TS-pooled was able to arrive at the correct conclusion (i. e., reject the null hypothesis of equal mean) $= 8$

- Number of times TS-**unpooled** performed **better** than TS-pooed (i.e., the p-value from TS-unpooled was **smaller** than that from TS-pooled)$= 0$

- Number of times TS-**pooled** performed **better** than TS-unpooed (i.e., the p-value from TS-pooled was **smaller** than that from TS-unpooled)$=$ 14

- Number of times TS-unpooled performed **as well** as TS-pooed (i.e., the p-value from TS-unpooled was **equal** to that from TS-pooled) = 1.

**Stage 2:**

| Population 1: $\mu_1 = 75$     $\sigma_1 = 15$ | Population 2: $\mu_2 = 65$     $\sigma_2 = 11$ |
|---|---|
| $\downarrow$ | $\downarrow$ |
| $n_1 = 50$ or $70$ | $n_2 = 60$ |

- Number of repeated samples = 15

- **Correct** Test Statistic: **TS-unpooled**

- Number of times TS-unpooled was able to arrive at the correct conclusion (i. e., reject the null hypothesis of equal mean) = 15

- Number of times TS-pooled was able to arrive at the correct conclusion (i. e., reject the null hypothesis of equal mean) = 15

- Number of times TS-**unpooled** performed **better** than TS-pooed (i.e., the p-value from TS-unpooled was **smaller** than that from TS-pooled) = 0

- Number of times TS-**pooled** performed **better** than TS-unpooed (i.e., the p-value from TS-pooled was **smaller** than that from TS-unpooled) = 5

- Number of times TS-unpooled performed **as well** as TS-pooed (i.e., the p-value from TS-unpooled was **equal** to that from TS-pooled) = 10.

**Stage 3:**

| Population 1: $\mu_1 = 75$     $\sigma_1 = 15$ | Population 2: $\mu_2 = 73$     $\sigma_2 = 11$ |
|---|---|
| $\downarrow$ | $\downarrow$ |
| $n_1 = 50$ or $70$ | $n_2 = 60$ |

- Number of repeated samples = 15

- **Correct** Test Statistic: **TS-unpooled**

- Number of times TS-unpooled was able to arrive at the correct conclusion (i. e., reject the null hypothesis of equal mean) = 1

- Number of times TS-pooled was able to arrive at the correct conclusion (i. e., reject the null hypothesis of equal mean) = 1

- Number of times TS-**unpooled** performed **better** than TS-pooed (i.e., the p-value from TS-unpooled was **smaller** than that from TS-pooled) = 8

- Number of times TS-**pooled** performed **better** than TS-unpooed (i.e., the p-value from TS-pooled was **smaller** than that from TS-unpooled) = 7

- Number of times TS-unpooled performed **as well** as TS-pooed (i.e., the p-value from TS-unpooled was **equal** to that from TS-pooled) = 0.

**Stage 4:**

| Population 1: $\mu_1 = 75$ | $\sigma_1 = 15$ | Population 2: $\mu_2 = 70$ | $\sigma_2 = 11$ |
|---|---|---|---|
| $\downarrow$ | | $\downarrow$ | |
| $n_1 = 50$ or $70$ | | $n_2 = 60$ | |

- Number of repeated samples = 15

- **Correct** Test Statistic: **TS-unpooled**

- Number of times TS-unpooled was able to arrive at the correct conclusion (i. e., reject the null hypothesis of equal mean) = 9

- Number of times TS-pooled was able to arrive at the correct conclusion (i. e., reject the null hypothesis of equal mean) = 10

- Number of times TS-**unpooled** performed **better** than TS-pooed (i.e., the p-value from TS-unpooled was **smaller** than that from TS-pooled) = 6

- Number of times TS-**pooled** performed **better** than TS-unpooed (i.e., the p-value from TS-pooled was **smaller** than that from TS-unpooled) = 7

- Number of times TS-unpooled performed **as well** as TS-pooed (i.e., the p-value from TS-unpooled was **equal** to that from TS-pooled) = 2

**Phase 2: The two populations have different means
but the same variance
(Stages 5 through 7 below)**

**Stage 5:**

| Population 1: $\mu_1 = 75$      $\sigma_1 = 15$ | Population 2: $\mu_2 = 70$      $\sigma_2 = 15$ |
|---|---|
| $\downarrow$ | $\downarrow$ |
| $n_1 = 50$ or $70$ | $n_2 = 75$ |

- Number of repeated samples $= 30$

- **Correct** Test Statistic: **TS-pooled**

- Number of times TS-unpooled was able to arrive at the correct conclusion (i. e., reject the null hypothesis of equal mean) $= 11$

- Number of times TS-pooled was able to arrive at the correct conclusion (i. e., reject the null hypothesis of equal mean) $= 11$

- Number of times TS-**unpooled** performed **better** than TS-pooed (i.e., the p-value from TS-unpooled was **smaller** than that from TS-pooled) $= 6$

- Number of times TS-**pooled** performed **better** than TS-unpooed (i.e., the p-value from TS-pooled was **smaller** than that from TS-unpooled) $= 14$

- Number of times TS-unpooled performed **as well** as TS-pooed (i.e., the p-value from TS-unpooled was **equal** to that from TS-pooled) $= 10$.

**Stage 6:**

| Population 1: $\mu_1 = 75$      $\sigma_1 = 15$ | Population 2: $\mu_2 = 65$      $\sigma_2 = 15$ |
|---|---|
| $\downarrow$ | $\downarrow$ |
| $n_1 = 50$ or $70$ | $n_2 = 75$ |

- Number of repeated samples $= 15$

- **Correct** Test Statistic: **TS-pooled**

- Number of times TS-unpooled was able to arrive at the correct conclusion (i. e., reject the null hypothesis of equal mean) $= 14$

- Number of times TS-pooled was able to arrive at the correct conclusion (i. e., reject the null hypothesis of equal mean) = 14

- Number of times TS-**unpooled** performed **better** than TS-pooed (i.e., the p-value from TS-unpooled was **smaller** than that from TS-pooled) = 1

- Number of times TS-**pooled** performed **better** than TS-unpooed (i.e., the p-value from TS-pooled was **smaller** than that from TS-unpooled) = 5

- Number of times TS-unpooled performed **as well** as TS-pooed (i.e., the p-value from TS-unpooled was **equal** to that from TS-pooled) = 9.

**Stage 7:**

| Population 1: $\mu_1 = 75$ $\sigma_1 = 15$ | Population 2: $\mu_2 = 73$ $\sigma_2 = 15$ |
|---|---|
| $\downarrow$ | $\downarrow$ |
| $n_1 = 50$ or $70$ | $n_2 = 75$ |

- Number of repeated samples = 15

- **Correct** Test Statistic: **TS-pooled**

- Number of times TS-unpooled was able to arrive at the correct conclusion (i. e., reject the null hypothesis of equal mean) = 0

- Number of times TS-pooled was able to arrive at the correct conclusion (i. e., reject the null hypothesis of equal mean) = 0

- Number of times TS-**unpooled** performed **better** than TS-pooed (i.e., the p-value from TS-unpooled was **smaller** than that from TS-pooled) = 8

- Number of times TS-**pooled** performed **better** than TS-unpooed (i.e., the p-value from TS-pooled was **smaller** than that from TS-unpooled) = 4

- Number of times TS-unpooled performed **as well** as TS-pooed (i.e., the p-value from TS-unpooled was **equal** to that from TS-pooled) = 3.

## Phase 3: The two populations have the same mean
## but different variance
## (Stages 8 and 9 below)

**Stage 8:**

| Population 1: $\mu_1 = 75$ | $\sigma_1 = 15$ | Population 2: $\mu_2 = 75$ | $\sigma_2 = 11$ |
|---|---|---|---|
| $\downarrow$ | | $\downarrow$ | |
| $n_1 = 50$ or $70$ | | $n_2 = 65$ | |

- Number of repeated samples $= 45$

- **Correct** Test Statistic: **TS-unpooled**

- Number of times TS-unpooled was able to arrive at the correct conclusion (i. e., fail to reject the null hypothesis of equal mean) $= 44$

- Number of times TS-pooled was able to arrive at the correct conclusion (i. e., fail to reject the null hypothesis of equal mean) $= 44$

- Number of times TS-**unpooled** performed **better** than TS-pooed (i.e., the p-value from TS-unpooled was **bigger** than that from TS-pooled) $= 28$

- Number of times TS-**pooled** performed **better** than TS-unpooed (i.e., the p-value from TS-pooled was **bigger** than that from TS-unpooled) $= 15$

- Number of times TS-unpooled performed **as well** as TS-pooed (i.e., the p-value from TS-unpooled was **equal** to that from TS-pooled) $= 2$.

**Stage 9:**

| Population 1: $\mu_1 = 75$ | $\sigma_1 = 15$ | Population 2: $\mu_2 = 75$ | $\sigma_2 = 7$ |
|---|---|---|---|
| $\downarrow$ | | $\downarrow$ | |
| $n_1 = 50$ or $70$ | | $n_2 = 65$ | |

- Number of repeated samples $= 15$

- **Correct** Test Statistic: **TS-unpooled**

- Number of times TS-unpooled was able to arrive at the correct conclusion (i. e., fail to reject the null hypothesis of equal mean) $= 14$

- Number of times TS-pooled was able to arrive at the correct conclusion (i. e., fail to reject the null hypothesis of equal mean) $= 14$

- Number of times TS-**unpooled** performed **better** than TS-pooed (i.e., the p-value from TS-unpooled was **bigger** than that from TS-pooled) $= 14$

- Number of times TS-**pooled** performed **better** than TS-unpooed (i.e., the p-value from TS-pooled was **bigger** than that from TS-unpooled) $= 0$

- Number of times TS-unpooled performed **as well** as TS-pooed (i.e., the p-value from TS-unpooled was **equal** to that from TS-pooled) $= 1$

## Phase 4: The two populations have the same mean and the same variance (Stage 10 below)

**Stage 10:**

| Population 1: $\mu_1 = 75$ | $\sigma_1 = 15$ | Population 2: $\mu_2 = 75$ | $\sigma_2 = 15$ |
|---|---|---|---|
| $\downarrow$ | | $\downarrow$ | |
| $n_1 = 50$ | | $n_2 = 70$ | |

- Number of repeated samples $= 60$

- **Correct** Test Statistic: **TS-pooled**

- Number of times TS-unpooled was able to arrive at the correct conclusion (i. e., fail to reject the null hypothesis of equal mean) $= 54$

- Number of times TS-pooled was able to arrive at the correct conclusion (i. e., fail to reject the null hypothesis of equal mean) $= 54$

- Number of times TS-**unpooled** performed **better** than TS-pooed (i.e., the p-value from TS-unpooled was **bigger** than that from TS-pooled) $= 20$

- Number of times TS-**pooled** performed **better** than TS-unpooed (i.e., the p-value from TS-pooled was **bigger** than that from TS-unpooled) $= 37$

- Number of times TS-unpooled performed **as well** as TS-pooed (i.e., the p-value from TS-unpooled was **equal** to that from TS-pooled) = 3

## 4. Conclusion

In this paper we use a simulation study to demonstrate that although both the pooled t-test and the unpooled t-test underperform at times, the overall performance of the pooled t-test is significantly supirior to that of the unpooled t-test. When the performances of the pooled t-test and the unpooled t-test were compared to each other, we founde the performance of the pooled t-test in the equal variance setting to be superior to the performance of the un-pooled t-test in the unequal-variance setting. Furthermore, when compared how the two test statistics perform when the underlying variance assumptions of their allocated settings were *not* met, we found the performance of the pooled t-test to be superior to the unpooled t-test.

## References

[1]  R. Peck, C. Olsen, J. Devore, *Introduction to Statistics and Data Analysis*, Brooks/Cole, USA (2009).

[2]  J. Devore, *Probability and Statistics for Engineering and the Sciences*, Brooks/Cole, USA (2008).

[3]  R. Johnson, G. Bhattacharyya, *Statistics-Principles and Methods*, Wiley, USA (2010).