

STAT 778: Midterm Exam

Tom Wallace

March 24, 2018

Preface: Program Organization and Compilation

Source code is contained in `midterm.c`. The program requires the GNU Scientific Library (GSL), an open-source numerical library. It can be obtained from www.gnu.org/software/gsl; or, it can be installed from most standard Linux package managers. An example command to achieve the latter is:

```
sudo apt-get install gsl-bin libgsl-dev
```

Compilation of `midterm.c` is best achieved in two steps. First, use the below command to compile the program but not link it. You may need to change the argument passed to the `-I` flag to wherever the `gsl` header files live on your computer.

```
gcc -I/usr/include -c midterm.c
```

This command should create an object file `midterm.o`. Link this object file to relevant libraries with the following command. You may need to change the argument passed to the `-L` flag to wherever `libgsl` lives on your computer.

```
gcc -L/usr/lib midterm.o -o midterm -lgsl -lgslcblas -lm
```

Once successfully compiled, the program can be executed. It does not require any arguments. Output is comma-separated text printed to `stdout`. You likely want to pipe this output to a text file, as per the following command:

```
./midterm > output.csv
```

Introduction

This study seeks to compare the performance of the two-sample t-test and the Wilcoxon rank-sum test. It uses simulation to do so. Data is randomly generated under different scenarios. For each scenario, the two methods are used to test the null hypothesis of no difference of means against the simple alternate hypothesis. The goal is to ascertain which method performs better by various criteria.

The remainder of this document is organized into two sections. The **Methods** section provides more detail on how the two methods were implemented and how their performance was compared. The **Simulation Study** section presents output data and results.

Methods

Tests for Difference of Means

The study compares the performance of two different tests of means. The first is **Welch's t-test**. This test assumes that the two populations are independent (i.e. unpaired), that they have normal distributions, and that they may have unequal variances. The test statistic t is calculated as:

$$t = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} \quad (1)$$

with \bar{X}_i , s_i^2 , and n_i denoting the sample mean, sample variance, and sample size of group i . The degrees of freedom for the t test statistic are calculated by the Welch-Satterthwaite equation:

$$df \approx \frac{\left(\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}\right)^2}{\frac{s_1^4}{n_1^2(n_1-1)} + \frac{s_2^4}{n_2^2(n_2-1)}} \quad (2)$$

The second is the **Wilcoxon rank-sum test**. This test makes no parametric assumptions nor any assumptions regarding common variance. Observations from the two groups are pooled, and then ranked in ascending order. The sum of ranks is taken for a group (which does not matter). The u statistic is given by:

$$u_1 = R_1 - \frac{n_1(n_1 + 1)}{2} \quad (3)$$

where R_1 is the sum of ranks of group 1, and n_1 is the sample size of group 1. This study uses the normal approximation for groups with $n_i \geq 25$:

$$z = \frac{u - \mu_u}{\sigma_u} \quad (4)$$

where $\mu_u = \frac{n_1 n_2}{2}$ and $\sigma_u = \sqrt{\frac{n_1 n_2 (n_1 + n_2 + 1)}{12}}$

Hypothesis Testing

The null hypothesis is no difference in group means, with a two-sided alternate hypothesis:

$$H_0 : \mu_1 = \mu_2$$

$$H_1 : \mu_1 \neq \mu_2$$

A significance level of $\alpha = 0.05$ is used.

Measures of Performance

The two tests are compared using the following measures, all of which will be explored using simulation:

- Type I error: at what rate does the test incorrectly reject the null?
- Power: at what rate does the test correctly reject the null (for various alternate parameter values)?
- Robustness to sample size: how sensitive is the test (in term of type I error and power) to sample size?
- Parametric robustness: how sensitive is the test (in terms of type I error and power) to different distributions?
- Robustness to outliers: how sensitive is the test (in terms of type I error and power) to outliers?

Simulation Study

Approach

The basic approach was to generate simulated data for two groups according to some specification; apply the t-test and Wilcoxon rank-sum test to the simulated data; and assess the type I error rate and power of each test. Different specifications were used, with variation in group size, distribution (including presence of outliers), and true difference in means. Each specification was simulated 1000 times. The tested specifications included various combinations of the values presented in Table 1.

Table 1: Specification Values

Distribution	n	True difference
Normal	25	None
Normal, 2% chance of outlier ($\mu_{\text{outlier}} = 10\mu$)	50	$\mu_2 = 1.1\mu_1$
Exponential	100	$\mu_2 = 1.25\mu_1$

Results

Quantitative results are presented in Table 2.

Comments Regarding Course

Overall, I have a positive impression of the class thus far. My C programming skills have gotten much sharper. I particularly appreciate that the instructor focuses on more on students learning useful skills than on grades. The purpose of graduate school is to prepare students to do professional research, and so I feel that grades are relatively superfluous. They can be a useful barometer of whether students are learning the skills needed for research, but the relationship is not that strong. In some other statistics courses, the pressure to do well on frequent, difficult graded assignments actually hurt my progress as a statistician. I got excellent grades, but did so by focusing more on learning to crank out correct answers than on truly understanding the material. This is a long way of saying: I really appreciate the different approach of this class. I put in just as much work, but am free to do so in a way that aids my long-term progress.

I hope that in the second half of the course we spend more time on algorithms. Techniques such as jackknife, bootstrap, EM, and the like are fundamental to modern statistics. Most of our work in the first half of the class has been on learning C, not learning statistical algorithms (we have been implementing very basic procedures).

Table 2: Simulation Results

Distribution	n	μ_1	σ_1^2	μ_2	σ_2^2	$\bar{\alpha}_t$	$\bar{\beta}_t$	$1 - \bar{\beta}_t$	$\bar{\alpha}_u$	$\bar{\beta}_u$	$1 - \bar{\beta}_u$
Normal	25	1.0	1.0	1.00	1.0	foo	bar	foo	bar	foo	bar
				1.10		foo	bar	foo	bar	foo	bar
				1.25		foo	bar	foo	bar	foo	bar
Normal	50	1.0	1.0	1.00	1.0	foo	bar	foo	bar	foo	bar
				1.10		foo	bar	foo	bar	foo	bar
				1.25		foo	bar	foo	bar	foo	bar
Normal	100	1.0	1.0	1.00	1.0	foo	bar	foo	bar	foo	bar
				1.10		foo	bar	foo	bar	foo	bar
				1.25		foo	bar	foo	bar	foo	bar
Normal (2% outlier)	25	1.0	1.0	1.00	1.0	foo	bar	foo	bar	foo	bar
				1.10		foo	bar	foo	bar	foo	bar
				1.25		foo	bar	foo	bar	foo	bar
Normal (2% outlier)	50	1.0	1.0	1.00	1.0	foo	bar	foo	bar	foo	bar
				1.10		foo	bar	foo	bar	foo	bar
				1.25		foo	bar	foo	bar	foo	bar
Normal (2% outlier)	100	1.0	1.0	1.00	1.0	foo	bar	foo	bar	foo	bar
				1.10		foo	bar	foo	bar	foo	bar
				1.25		foo	bar	foo	bar	foo	bar
Exponential	25	1.0	1.0	1.00	1.00	foo	bar	foo	bar	foo	bar
				1.10	1.10	foo	bar	foo	bar	foo	bar
				1.25	1.25	foo	bar	foo	bar	foo	bar
Exponential	50	1.0	1.0	1.00	1.00	foo	bar	foo	bar	foo	bar
				1.10	1.10	foo	bar	foo	bar	foo	bar
				1.25	1.25	foo	bar	foo	bar	foo	bar
Exponential	100	1.0	1.0	1.00	1.00	foo	bar	foo	bar	foo	bar
				1.10	1.10	foo	bar	foo	bar	foo	bar
				1.25	1.25	foo	bar	foo	bar	foo	bar

Each specification simulated 1000 times

$\bar{\alpha}_t$ = empirical Type I error probability, t-test

$\bar{\alpha}_u$ = empirical Type I error probability, Wilcoxon rank-sum test

$\bar{\beta}_t$ = empirical Type II error probability, t-test

$\bar{\beta}_u$ = empirical Type II error probability, Wilcoxon rank-sum test