# STAT 401A - Statistical Methods for Research Workers
## Modeling assumptions

Jarad Niemi (Dr. J)

Iowa State University

last updated: September 11, 2014

# Normality assumptions

In the paired t-test, we assume

$$D_i \overset{iid}{\sim} N(\mu, \sigma^2).$$

In the two-sample t-test, we assume

$$Y_{ij} \overset{ind}{\sim} N(\mu_j, \sigma^2).$$
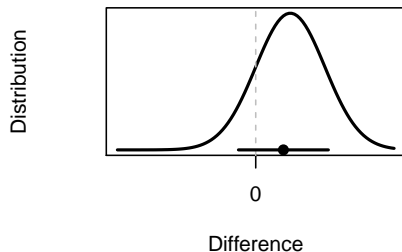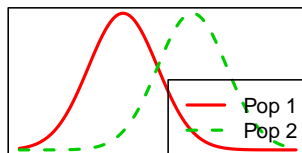
# Normality assumptions

In the paired t-test, we assume

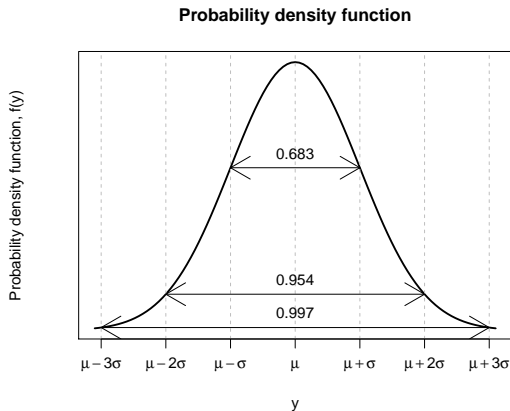$$D_i \stackrel{iid}{\sim} N(\mu, \sigma^2).$$

In the two-sample t-test, we assume

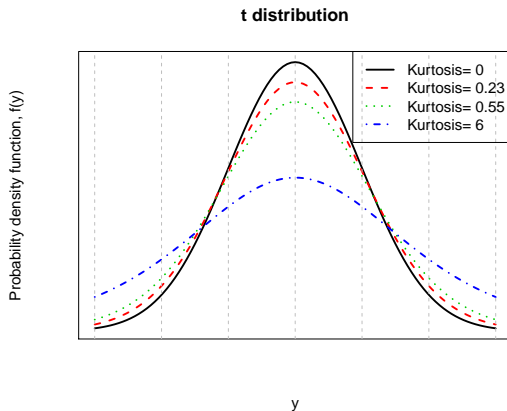$$Y_{ij} \stackrel{ind}{\sim} N(\mu_j, \sigma^2).$$

Key features of the normal distribution assumption:

- Centered at the mean (expectation) $\mu$
- Standard deviation describes the spread
- Symmetric around $\mu$ (no skewness)
- Non-heavy tails, i.e. outliers are rare (no kurtosis)

# Normality assumptions



**Probability density function**

# Kurtosis (heavy-tailedness)



**t distribution**

# Kurtosis (heavy-tailedness)



**Probability density function**

# Kurtosis (heavy-tailedness)

# Kurtosis (heavy-tailedness)

# Skewness

**Log–normal distribution**

# Samples from skewed distributions

# Robustness

### Definition
A statistical procedure is robust to departures from a particular assumption if it is valid even when the assumption is not met.

**Remark** If a 95% confidence interval is robust to departures from a particular assumption, the confidence interval should cover the true value about 95% of the time.

# Robustness to skewness and kurtosis

Percentage of 95% confidence intervals that cover the true difference in means in an equal-sample two-sample t-test with non-normal populations (where the distributions are the same other than their means).
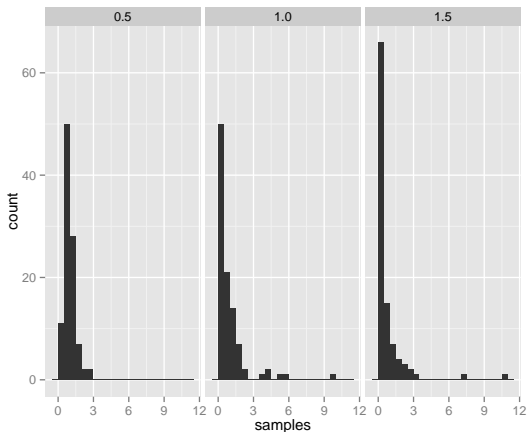
| sample size | strongly skewed | moderately skewed | mildly skewed | heavy-tailed | short-tailed |
|---|---|---|---|---|---|
| 5 | 95.5 | 95.4 | 95.2 | 98.3 | 94.5 |
| 10 | 95.5 | 95.4 | 95.2 | 98.3 | 94.6 |
| 25 | 95.3 | 95.3 | 95.1 | 98.2 | 94.9 |
| 50 | 95.1 | 95.3 | 95.1 | 98.1 | 95.2 |
| 100 | 94.8 | 95.3 | 95.0 | 98.0 | 95.6 |

# Differences in variances

**Normal distribution**
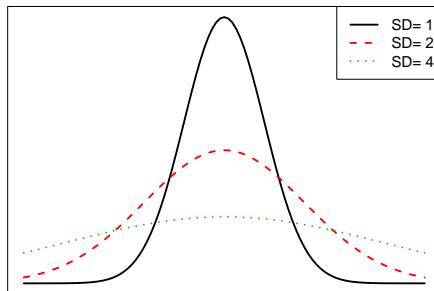
# Differences in variances

# Robustness to differences in variances

Percentage of 95% confidence intervals that cover the true difference in means in an equal-sample two-sample t-test ($r = \sigma_1/\sigma_2$).

| n1 | n2 | r=1/4 | r=1/2 | r=1 | r=2 | r=4 |
|----|----|-------|-------|-----|-----|-----|
| 10 | 10 | 95.2 | 94.2 | 94.7 | 95.2 | 94.5 |
| 10 | 20 | 83.0 | 89.3 | 94.4 | 98.7 | 99.1 |
| 10 | 40 | 71.0 | 82.6 | 95.2 | 99.5 | 99.9 |
| 100 | 100 | 94.8 | 96.2 | 95.4 | 95.3 | 95.1 |
| 100 | 200 | 86.5 | 88.3 | 94.8 | 98.8 | 99.4 |
| 100 | 400 | 71.6 | 81.5 | 95.0 | 99.5 | 99.9 |

# Outliers

### Definition
A statistical procedure is resistant if it does not change very much when a small part of the data changes, perhaps drastically.

Identify outliers:

1. If recording errors, fix.
2. If outlier comes from a different population, remove and report.
3. If results are the same with and without outliers, report with outliers.
4. If results are different, use resistant analysis or report both analyses.

# Common ways for independence to be violated

- Cluster effect
  - e.g. pigs in a pen
- Correlation effect
  - e.g. measurements in time with drifting scale
- Spatial effect
  - e.g. corn yield plots (drainage)

# Common transformations for data

From: http://en.wikipedia.org/wiki/Data_transformation_(statistics)

### Definition

In statistics, data transformation refers to the application of a deterministic mathematical function to each point in a data set  that is, each data point $y_i$ is replaced with the transformed value $z_i = f(y_i)$, where $f$ is a function.

The most common transformations to

- If $y \in (0, 1)$, then $f(y) = sin^{-1}(\sqrt{y})$.
- If $y$ is a count, then $f(y) = \sqrt{y}$.
- If $y$ is positive and right-skewed, then $f(y) = log(y)$, the *natural logarithm* of $y$.

**Remark** Since $log(0) = -\infty$, the logarithm cannot be used directly when some $y_i$ are zero. In these cases, use $log(y + c)$ where $c$ is something small relative to your data, e.g. half of the minimum non-zero value.

# Log transformation

If $z_{ij} = log(y_{ij})$ and we run a two-sample t-test on the z's, then we assume

$$Z_{ij} \stackrel{ind}{\sim} N(\mu_j, \sigma^2)$$

and we calculate the quantity $\overline{Z}_2 - \overline{Z}_1$ and $\exp\left(\overline{Z}_2 - \overline{Z}_1\right) = e^{\overline{Z}_2 - \overline{Z}_1}$ estimates

$$\frac{\text{Median of population 2}}{\text{Median of population 1}}$$

on the original scale or, equivalently, it estimates the multiplicative effect of moving from population 1 to population 2.

# Log transformation interpretation

If we have a randomized experiment:

**Remark** It is estimated that the response of an experimental unit to treatment 2 will be $\exp\left(\overline{Z}_2 - \overline{Z}_1\right)$ times as large as its response to treatment 1.

If we have an observational study:

**Remark** It is estimated that the median for population 2 is $\exp\left(\overline{Z}_2 - \overline{Z}_1\right)$ times as large as the median for population 1.

# Confidence intervals with log transformation

If $z_{ij} = log(y_{ij})$ and we assume

$$Z_{ij} \overset{ind}{\sim} N(\mu_j, \sigma^2),$$

then a $100(1 - \alpha)\%$ two-sided confidence interval for $\mu_2 - \mu_1$ is

$$(L, U) = \overline{Z}_2 - \overline{Z}_1 \pm t_{n_1+n_2-2}(1 - \alpha/2) SE\left(\overline{Z}_2 - \overline{Z}_1\right).$$
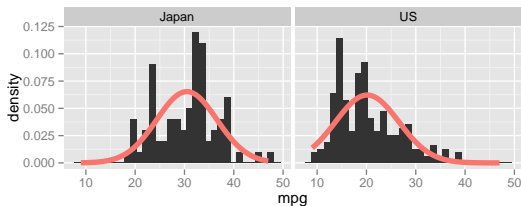
A $100(1 - \alpha)\%$ confidence interval for

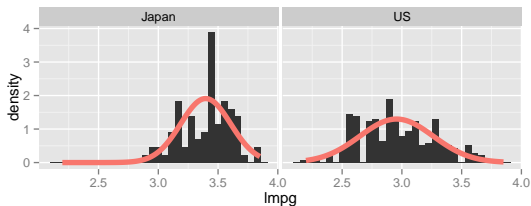$$\frac{\text{Median of population 2}}{\text{Median of population 1}}$$

is $(e^L, e^U)$.

# Miles per gallon data

Untransformed:



Logged:

# Miles per gallon data

Untransformed:



Logged:

# Equal variances?

We might also be concerned about the assumption of equal variances.

Untransformed:

|   | country | n | mean | sd |
|---|---------|-----|-------|-------|
| 1 | Japan | 79 | 30.48 | 30.48 |
| 2 | US | 249 | 20.14 | 20.14 |

the ratio of the standard deviations is around 1.5 and there are 3 times as many observations in the US.

Logged:

|   | country | n | mean | sd |
|---|---------|-----|------|------|
| 1 | Japan | 79 | 3.40 | 3.40 |
| 2 | US | 249 | 2.96 | 2.96 |

Now the ratio of standard deviations is only 1.15.

# Using R for t-test using logarithms

```
t = t.test(log(mpg)~country, d, var.equal=TRUE)
exp(t$estimate)

mean in group Japan    mean in group US
            29.85                  19.21


exp(-diff(t$estimate)) # I had to put in the negative sign

mean in group US
         1.554


exp(t$conf.int)

[1] 1.445 1.672
attr(,"conf.level")
[1] 0.95
```

# SAS code for t-test using logarithms

```
DATA mpg;
   INFILE 'mpg.csv' DELIMITER=',' FIRSTOBS=2;
INPUT mpg country $;

PROC TTEST DATA=mpg TEST=ratio;
CLASS country;
VAR mpg;
run;
```

# SAS output for t-test using logarithms

The TTEST Procedure

Variable: mpg

| country | N | Geometric Mean | Coefficient of Variation | Minimum | Maximum |
|---------|---|----------------|--------------------------|---------|---------|
| Japan | 79 | 29.8525 | 0.2111 | 18.0000 | 47.0000 |
| US | 249 | 19.2051 | 0.3147 | 9.0000 | 39.0000 |
| Ratio (1/2) | | 1.5544 | 0.2928 | | |

| country | Method | Geometric Mean | 95% CL Mean | | Coefficient of Variation | 95% CL CV | |
|---------|--------|----------------|-------------|---|--------------------------|-----------|---|
| Japan | | 29.8525 | 28.4887 | 31.2817 | 0.2111 | 0.1820 | 0.2514 |
| US | | 19.2051 | 18.4825 | 19.9560 | 0.3147 | 0.2882 | 0.3467 |
| Ratio (1/2) | Pooled | 1.5544 | 1.4452 | 1.6719 | 0.2928 | 0.2712 | 0.3183 |
| Ratio (1/2) | Satterthwaite | 1.5544 | 1.4636 | 1.6508 | | | |

| Method | Coefficients of Variation | DF | t Value | Pr > |t| |
|--------|---------------------------|------|---------|----------|
| Pooled | Equal | 326 | 11.91 | <.0001 |
| Satterthwaite | Unequal | 193.33 | 14.46 | <.0001 |

Equality of Variances

| Method | Num DF | Den DF | F Value | Pr > F |
|--------|--------|--------|---------|--------|
| Folded F | 248 | 78 | 2.17 | 0.0001 |

# Conclusion

Japanese median miles per gallon is 1.55 [95% CI (1.46,1.65)] times as large as US median miles per gallon.