

STAT 401A - Statistical Methods for Research Workers

Modeling assumptions

Jarad Niemi (Dr. J)

Iowa State University

last updated: September 14, 2014

Normality assumptions

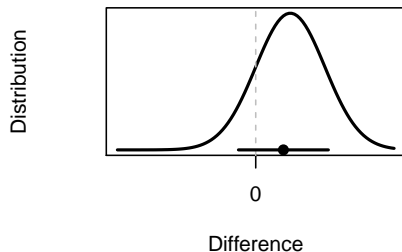
In the paired t-test, we assume

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

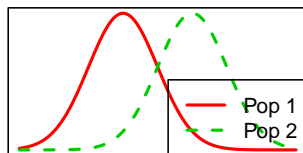
In the two-sample t-test, we assume

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

Paired t-test



Two-sample t-test



Normality assumptions

In the paired t-test, we assume

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

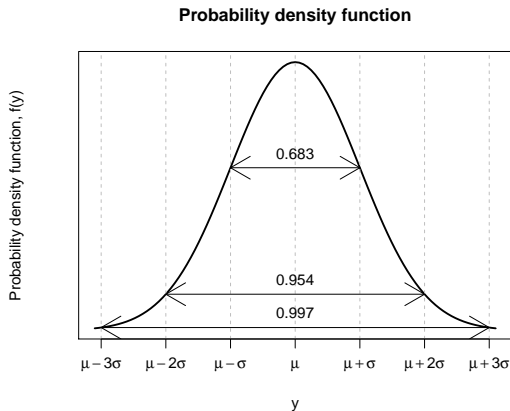
In the two-sample t-test, we assume

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

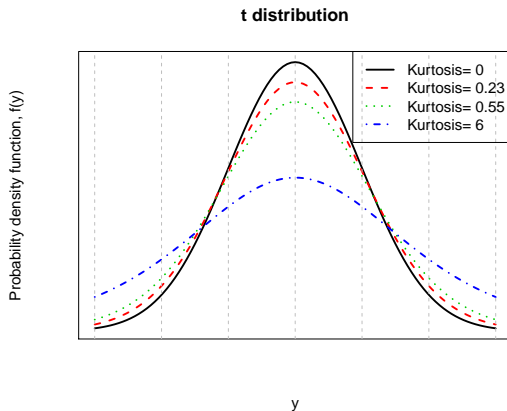
Key features of the normal distribution assumption:

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

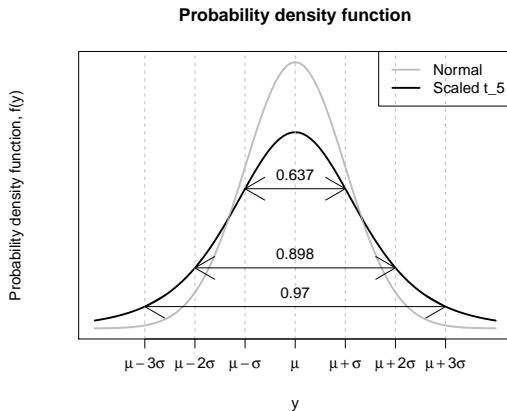
Normality assumptions



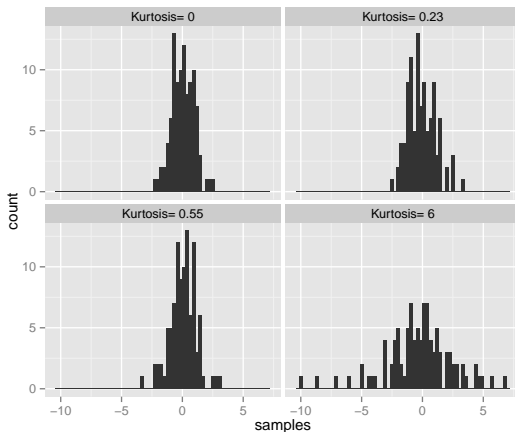
Kurtosis (heavy-tailedness)



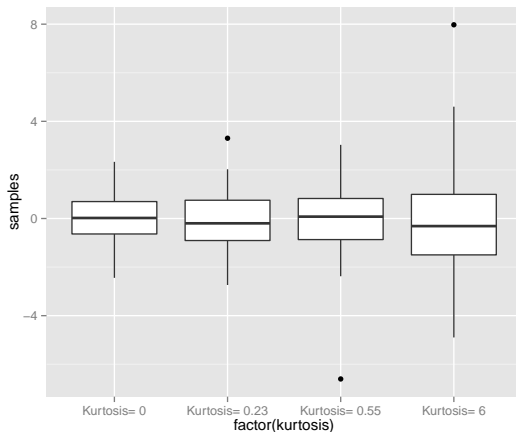
Kurtosis (heavy-tailedness)



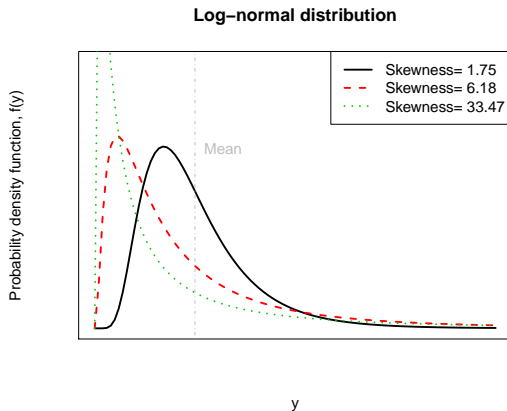
Kurtosis (heavy-tailedness)



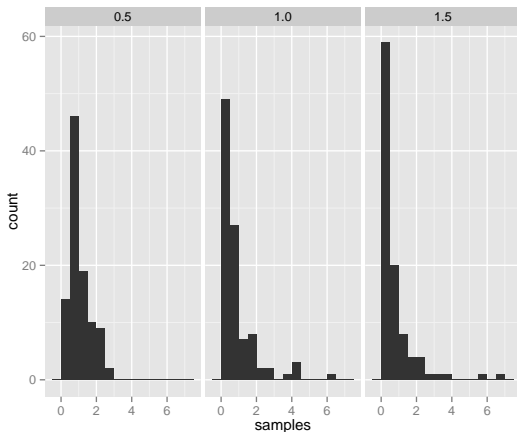
Kurtosis (heavy-tailedness)



Skewness



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.

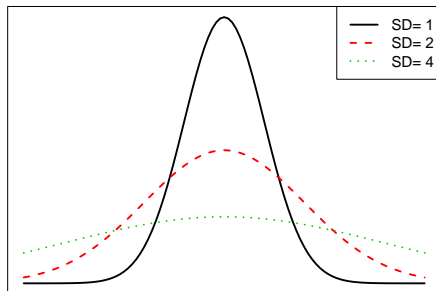
Robustness to skewness and kurtosis

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).

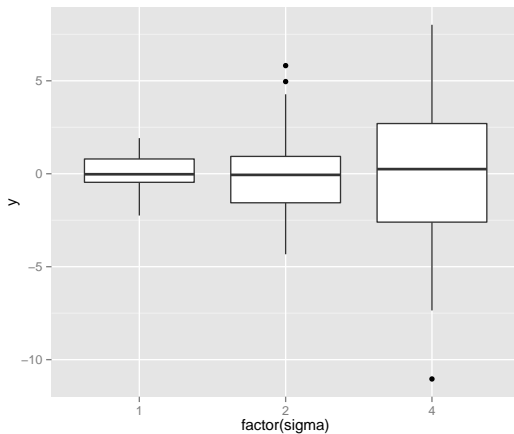
sample size	strongly skewed	moderately skewed	mildly skewed	heavy-tailed	short-tailed
5	95.5	95.4	95.2	98.3	94.5
10	95.5	95.4	95.2	98.3	94.6
25	95.3	95.3	95.1	98.2	94.9
50	95.1	95.3	95.1	98.1	95.2
100	94.8	95.3	95.0	98.0	95.6

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_\(statistics\)](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 are

- If y is a proportion, 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

Consider two-sample data and let $z_{ij} = \log(y_{ij})$. Now, run a two-sample t-test on the z 's. Then we assume

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

and the quantity $\bar{Z}_2 - \bar{Z}_1$ estimates the “difference in population means on the (natural) log scale”. The quantity $\exp(\bar{Z}_2 - \bar{Z}_1) = e^{\bar{Z}_2 - \bar{Z}_1}$ 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 $\exp(\bar{Z}_2 - \bar{Z}_1)$ 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 $\exp(\bar{Z}_2 - \bar{Z}_1)$ times as large as the median for population 1.

Confidence intervals with log transformation

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

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

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

$$(L, U) = \bar{Z}_2 - \bar{Z}_1 \pm t_{n_1+n_2-2}(1 - \alpha/2)SE(\bar{Z}_2 - \bar{Z}_1).$$

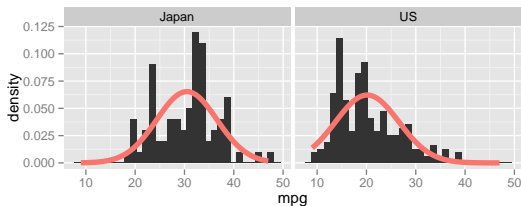
A $100(1 - \alpha)\%$ confidence interval for

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

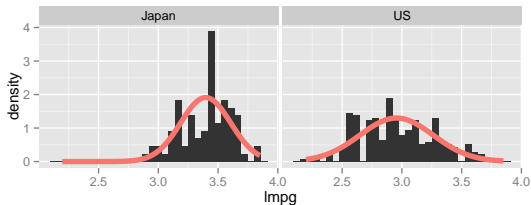
is (e^L, e^U) .

Miles per gallon data

Untransformed:

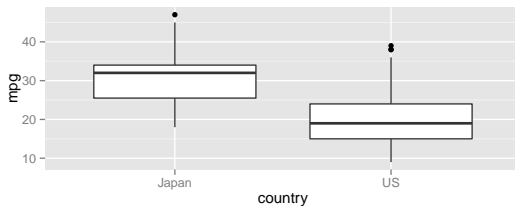


Logged:

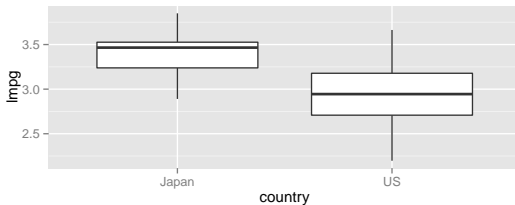


Miles per gallon data

Untransformed:



Logged:



Equal variances?

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

Untransformed:

country	n	mean	sd
Japan	79	30.48	6.11
US	249	20.14	6.41

the ratio of sample standard deviations is around 1.05 and there are 3 times as many observations in the US.

Logged:

country	n	mean	sd
Japan	79	3.40	0.21
US	249	2.96	0.31

Now the ratio of standard deviations is only 1.5 which argues for not using the logarithm.

95% two-sample CI for the ratio by hand

country	n	mean	sd
Japan	79	3.40	0.21
US	249	2.96	0.31

Choose group 2 to be Japan and group 1 to be the US:

$$\begin{aligned}
 \alpha &= 0.05 \\
 n_1 + n_2 - 2 &= 249 + 79 - 2 = 326 \\
 t_{n_1+n_2-2}(1 - \alpha/2) &= t_{326}(0.975) = 1.96 \\
 \bar{Z}_2 - \bar{Z}_1 &= 3.40 - 2.96 = 0.44 \\
 s_p &= \sqrt{\frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{n_1+n_2-2}} = \sqrt{\frac{(249-1)0.31^2 + (79-1)0.21^2}{249+79-2}} = 0.29 \\
 SE(\bar{Z}_2 - \bar{Z}_1) &= s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}} = 0.29 \sqrt{\frac{1}{249} + \frac{1}{79}} = 0.037
 \end{aligned}$$

Thus a 95% two-sided confidence interval for the difference (on the log scale) is

$$\begin{aligned}
 (L, U) &= \bar{Z}_2 - \bar{Z}_1 \pm t_{n_1+n_2-2}(1 - \alpha/2)SE(\bar{Z}_2 - \bar{Z}_1) \\
 &= 0.44 \pm 1.96 \times 0.037 \\
 &= (0.37, 0.51)
 \end{aligned}$$

and a 95% two-sided confidence interval for the ratio (on the original scale) is

$$(e^L, e^U) = (e^{0.37}, e^{0.51}) = (1.45, 1.67)$$

Using R for t-test using logarithms

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

```
t$estimate # On log scale
```

mean in group Japan	mean in group US
3.396	2.955

```
exp(t$estimate) # On original scale
```

mean in group Japan	mean in group US
29.85	19.21

```
exp(t$estimate[1]-t$estimate[2]) # Ratio of medians (Japan/US)
```

mean in group Japan
1.554

```
exp(t$conf.int) # Confidence interval for ratio of medians
```

```
[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.

OR

Japanese median miles per gallon is 55% [95% CI (46%,65%)] larger than US median miles per gallon.