

R refresher

January 24, 2022

R versus other languages

- Weird assignment operator <-
- cAsE sENsiTiVe
- Indexing starts at 1 instead of 0 (first element in a vector is element 1)
- Always more than one way to do things (sometimes many more)
- The R Inferno
 - https://www.burns-stat.com/pages/Tutor/R_inferno.pdf

Basics

- Assignment `x <- 5`
- Arithmetic operations `3 + 5, x + 3`
- Vectors `y <- c(1, 4, 5, 9)`
- Lists `z <- list(1, "wrong", c(3, 5, 7))`
- data types - "double", integer, logical, character, complex
- data structures – vector, list, matrix, data frame, factors, tables, "tibbles", ...
- `my_data_frame$variable1[1:4]` yuck
- functions!

Functions

```
add_five <- function(x) {x + 5}
```

```
> add_five(30)
```

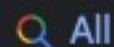
```
[1] 35
```

- Powerful and flexible
- "Functional programming"
- Work out how to do something once then wrap in a function

Figuring things out

- R help files are often not very helpful
- There is an enormous amount of helpful material online – if you have a problem, someone else has probably had that problem and figured out the answer
- Sometimes it is hard to determine whether you just can't do something or whether you haven't found the right answer yet
- Google is your friend "make square ggplot"

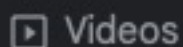
make square ggplot



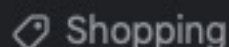
All



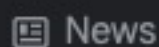
Images



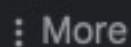
Videos



Shopping



News



More

Tools

About 355,000 results (0.68 seconds)

[https://stackoverflow.com > questions > force-ggplot2-s...](https://stackoverflow.com/questions/force-ggplot2-s...)

Force ggplot2 scatter plot to be square shaped - Stack Overflow

Mar 10, 2016 — I can force **ggplot2** scatter plot to be **square** shaped with the same x and y scaling using `xlim()` and `ylim()` , but it needs manual calculation of ...

5 answers · Top answer: If you want to make the distance scale points the same, then use co...

[How to fix the aspect ratio in ggplot? - Stack Overflow](#)

4 answers

Oct 6, 2013

[Equal coordinates and square aspect ratio with log scale ...](#)

1 answer

May 7, 2020

[Draw multiple squares with ggplot - Stack Overflow](#)

2 answers

Apr 9, 2013

[How do I make my facets perfectly square? - Stack ...](#)

4 answers

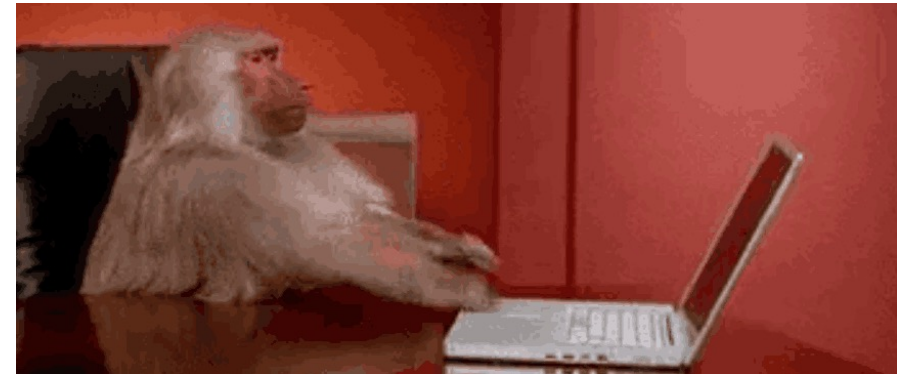
Jan 22, 2014

[More results from stackoverflow.com](#)

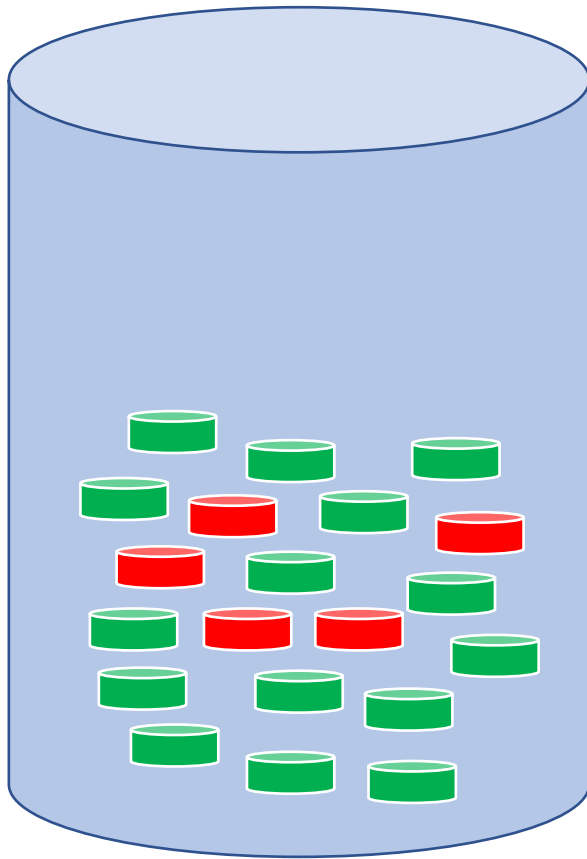
You've visited this page many times. Last visit: 12/22/21

Figuring things out

- You often just have to keep flailing at the keyboard to figure out how to do something or how to get something to work
- Walk away, come back, do something else for a little while
- Ask a friend if Google fails you
- You will make mistakes. This is inevitable. What you want to avoid is realizing you've made a mistake after the paper gets published.



Questions and discussion on "pre-test"



* probabilities sum to 1

$$\text{draw 1 chip: } P(\text{red}) = \frac{\text{number of red chips}}{\text{total chips}} = 0.25$$

$$P(\text{green}) = \frac{\text{number of green chips}}{\text{total chips}} = 0.75$$

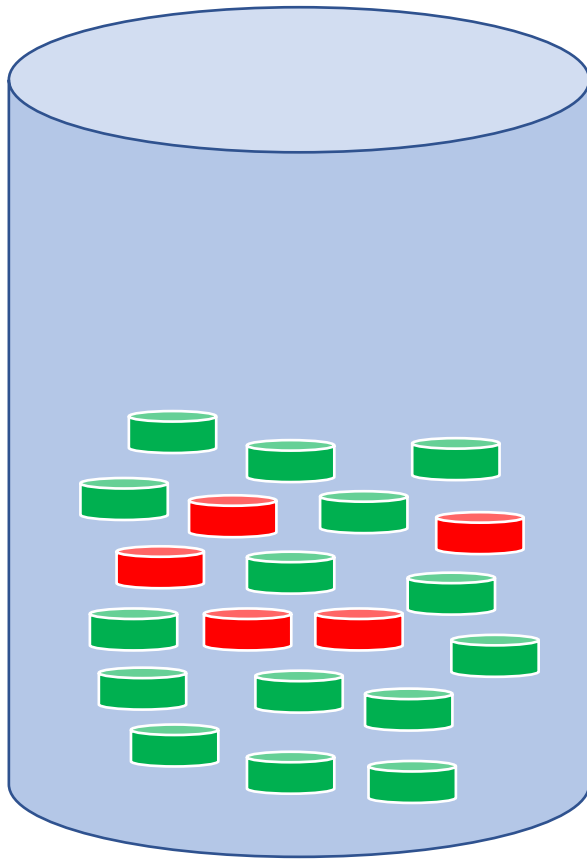
* sample with/without replacement

* $P(\text{all green}) + P(\text{at least 1 red}) = 1$

because one or the other must happen

$$P(\text{at least 1 red}) = 1 - P(\text{all green})$$

Questions and discussion on "pre-test"

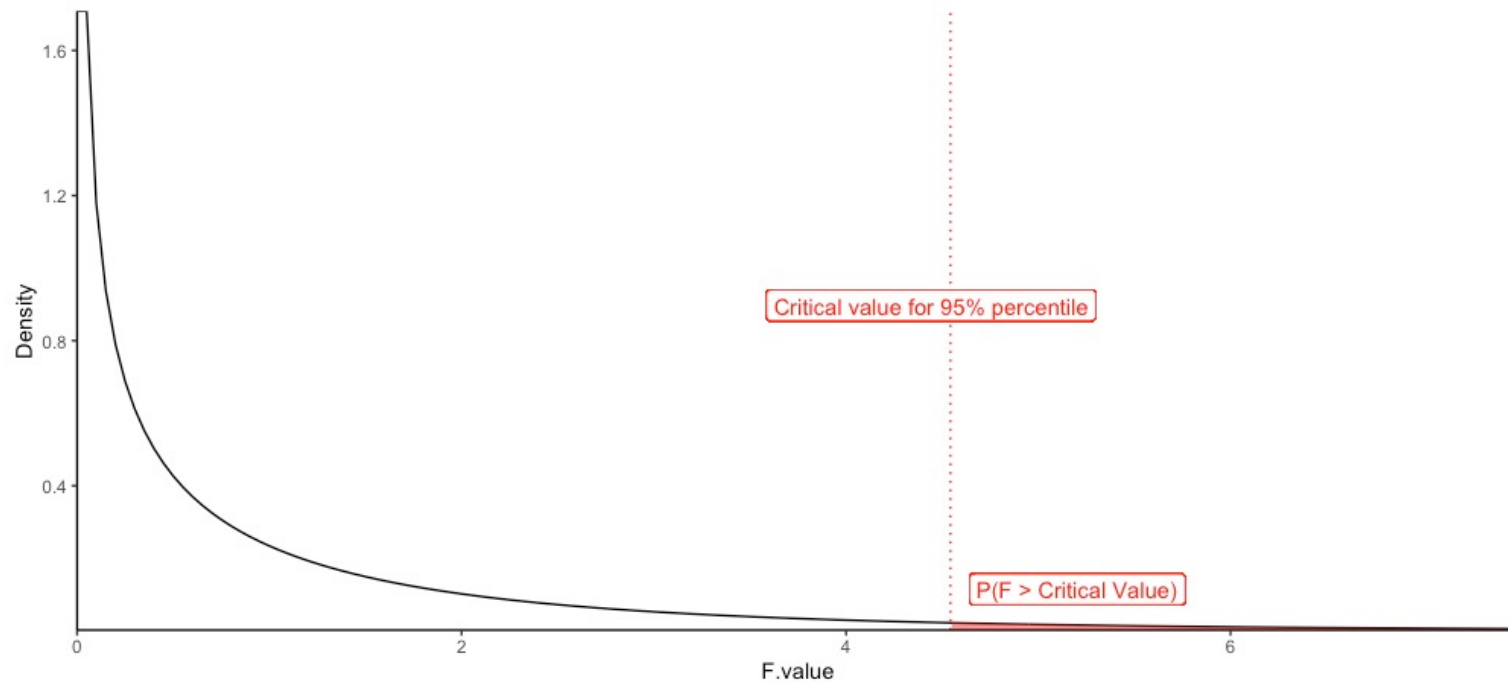


$$\begin{aligned} & * P(2 \text{ green}) = \\ & P(\text{green, green, red}) + \\ & P(\text{green, red, green}) + \\ & P(\text{red, green, green}) = \\ & 3 * (15/20) * (15/20) * (5/20) \end{aligned}$$

- * probabilities multiply
- * independent events add
- * order matters "combinatorics"

$$\binom{n}{r} = \frac{n!}{r!(n-r)!}$$

Questions and discussion on "pre-test"



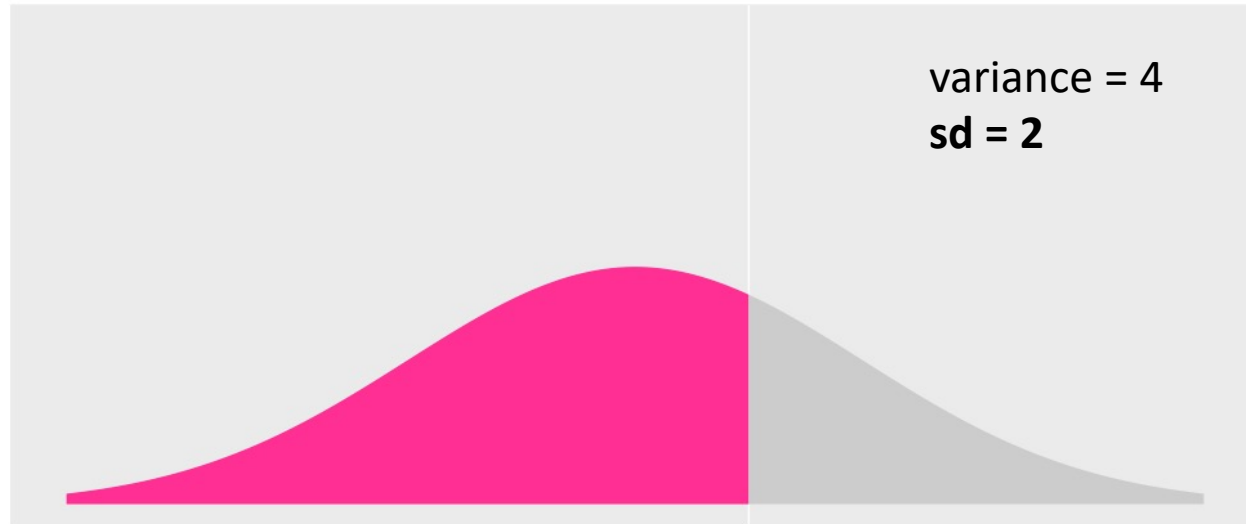
Idea of "critical value"

Getting areas from distribution (pf) versus getting value from curve (qf, df)

pnorm, qnorm, dnorm

Graphical representation – different graphics modules in R

Questions and discussion on "pre-test"



```
> pnorm(1, mean = 0, sd = 1, lower.tail = TRUE)
[1] 0.8413447
> pnorm(1, mean = 0, sd = 2)
[1] 0.6914625
```

Questions and discussion on "pre-test"

```
# imagine that these are test score data of 2 groups of participants, 10 in each group

group1 <- c(56, 46, 45, 42, 60, 45, 52, 59, 43, 55)
group2 <- c(85, 61, 57, 53, 64, 58, 67, 54, 76, 63)

# run a t-test comparing group1 and group2, ** assuming equal variances of the 2 groups **

t.test(group1, group2, var.equal = TRUE)

# What is the type 1 error associated with the t-test? (either copy from output or use an R command to display)

0.002518
t.test(group1, group2, var.equal = TRUE)$p.value # look at output of str(t.test(group1, group2, var.equal = TRUE))
```

- default for t.test is that variances are not equal
- attributes of object can be revealed by `str` ("structure") and accessed directly



?t.test at prompt

Usage

```
t.test(x, ...)
```

```
## Default S3 method:
t.test(x, y = NULL,
       alternative = c("two.sided", "less", "greater"),
       mu = 0, paired = FALSE, var.equal = FALSE,
       conf.level = 0.95, ...)
```

```
## S3 method for class 'formula'
t.test(formula, data, subset, na.action, ...)
```

FilesPlotsPackagesHelpViewer

←→🏠🖨️

🔍

Refresh Help Topic

R: Student's t-TestFind in Topic

t.test {stats}R Documentation

Student's t-Test

Description

Performs one and two sample t-tests on vectors of data.

Usage

```
t.test(x, ...)  
  
## Default S3 method:  
t.test(x, y = NULL,  
       alternative = c("two.sided", "less", "greater"),  
       mu = 0, paired = FALSE, var.equal = FALSE,  
       conf.level = 0.95, ...)  
  
## S3 method for class 'formula'  
t.test(formula, data, subset, na.action, ...)
```

Arguments

x	a (non-empty) numeric vector of data values.
y	an optional (non-empty) numeric vector of data values.
alternative	a character string specifying the alternative hypothesis, must be one of "two.sided" (default), "greater" or "less". You can specify just the initial letter.
mu	a number indicating the true value of the mean (or difference in means if you are performing a two sample test).
paired	a logical indicating whether you want a paired t-test.
var.equal	a logical variable indicating whether to treat the two variances as being equal. If TRUE then the pooled variance is used to estimate the variance otherwise the Welch (or Satterthwaite) approximation to the degrees of freedom is used.
conf.level	confidence level of the interval.
formula	a formula of the form lhs ~ rhs where lhs is a numeric variable giving the data values and rhs either 1 for a one-sample or paired test or a factor with two levels giving the corresponding groups. If lhs is of class "Pair" and rhs is 1, a paired test is done
data	an optional matrix or data frame (or similar: see <code>model.frame</code>) containing the variables in the formula <code>formula</code> . By default the variables are taken from <code>environment(formula)</code> .
subset	an optional vector specifying a subset of observations to be used.
na.action	a function which indicates what should happen when the data contain NAs. Defaults to <code>getOption("na.action")</code> .
...	further arguments to be passed to or from methods.

Details

`alternative = "greater"` is the alternative that `x` has a larger mean than `y`. For the one-sample case: that the mean is positive.

If `paired` is TRUE then both `x` and `y` must be specified and they must be the same length. Missing values are silently removed (in pairs if `paired` is TRUE). If `var.equal` is TRUE then the pooled estimate of the variance is used. By default, if `var.equal` is FALSE then the variance is estimated separately for both groups and the Welch modification to the degrees of freedom is used.

If the input data are effectively constant (compared to the larger of the two means) an error is generated.

Value

A list with class "htest" containing the following components:

statistic	the value of the t-statistic.
parameter	the degrees of freedom for the t-statistic.
p.value	the p-value for the test.

?t.test at prompt

Usage

```
t.test(x, ...)
```

```
## Default S3 method:
```

```
t.test(x, y = NULL,  
       alternative = c("two.sided", "less", "greater"),  
       mu = 0, paired = FALSE, var.equal = FALSE,  
       conf.level = 0.95, ...)
```

```
## S3 method for class 'formula'
```

```
t.test(formula, data, subset, na.action, ...)
```

describes default behavior or how it behaves if you give it a "formula" object



?t.test at prompt

Usage

```
t.test(x, ...)
```

```
## Default S3 method:
```

```
t.test(x, y = NULL,
       alternative = c("two.sided", "less", "greater"),
       mu = 0, paired = FALSE, var.equal = FALSE,
       conf.level = 0.95, ...)
```

```
## S3 method for class 'formula'
```

```
t.test(formula, data, subset, na.action, ...)
```

lists arguments the function expects. If there is an equals sign, it has a default value that will be used if it is not specified.


```
> t.test(group1)
```

One Sample t-test

```
data: group1
t = 23.172, df = 9, p-value = 2.47e-09
alternative hypothesis: true mean is not equal to 0
95 percent confidence interval:
 45.38939 55.21061
sample estimates:
mean of x
 50.3
```

```
> t.test(group1, group2)
```

Welch Two Sample t-test

```
data: group1 and group2
t = -3.5069, df = 15.894, p-value = 0.002946
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 -21.665073 -5.334927
sample estimates:
mean of x mean of y
 50.3      63.8
```

```
> t.test()
Error in t.test.default() : argument "x" is missing, with no default
```

?t.test at prompt

Usage

```
t.test(x, ...)
```

```
## Default S3 method:
```

```
t.test(x, y = NULL,
       alternative = c("two.sided", "less", "greater"),
       mu = 0, paired = FALSE, var.equal = FALSE,
       conf.level = 0.95, ...)
```

```
## S3 method for class 'formula'
```

```
t.test(formula, data, subset, na.action, ...)
```

lists arguments the function expects. If there is an equals sign, it has a default value that will be used if it is not specified.

x a (non-empty) numeric vector of data values.

y an optional (non-empty) numeric vector of data values.

alternative a character string specifying the alternative hypothesis, must be one of "two.sided" (default), "greater" or "less". You can specify just the initial letter.

statistic the value of the t-statistic.

parameter the degrees of freedom for the t-statistic.

p.value the p-value for the test.


```
> t.test(group1)
```

One Sample t-test

```
data: group1
t = 23.172, df = 9, p-value = 2.47e-09
alternative hypothesis: true mean is not equal to 0
95 percent confidence interval:
 45.38939 55.21061
sample estimates:
mean of x
 50.3
```

```
> t.test(group1, group2)
```

Welch Two Sample t-test

```
data: group1 and group2
t = -3.5069, df = 15.894, p-value = 0.002946
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 -21.665073 -5.334927
sample estimates:
mean of x mean of y
 50.3      63.8
```

```
> t.test()
Error in t.test.default() : argument "x" is missing, with no default
```

?t.test at prompt

Usage

```
t.test(x, ...)
```

```
## Default S3 method:
```

```
t.test(x, y = NULL,
       alternative = c("two.sided", "less", "greater"),
       mu = 0, paired = FALSE, var.equal = FALSE,
       conf.level = 0.95, ...)
```

```
## S3 method for class 'formula'
```

```
t.test(formula, data, subset, na.action, ...)
```

lists arguments the function expects. If there is an equals sign, it has a default value that will be used if it is not specified.

You have to give it at least one vector of values ("x")

```
> t.test(group1)
```

One Sample t-test

```
data: group1
t = 23.172, df = 9, p-value = 2.47e-09
alternative hypothesis: true mean is not equal to 0
95 percent confidence interval:
 45.38939 55.21061
sample estimates:
mean of x
 50.3
```

```
> t.test(group1, group2)
```

Welch Two Sample t-test

```
data: group1 and group2
t = -3.5069, df = 15.894, p-value = 0.002946
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 -21.665073 -5.334927
sample estimates:
mean of x mean of y
 50.3      63.8
```

```
> t.test()
Error in t.test.default() : argument "x" is missing, with no default
```

?t.test at prompt

Usage

```
t.test(x, ...)
```

```
## Default S3 method:
```

```
t.test(x, y = NULL,
       alternative = c("two.sided", "less", "greater"),
       mu = 0, paired = FALSE, var.equal = FALSE,
       conf.level = 0.95, ...)
```

```
## S3 method for class 'formula'
```

```
t.test(formula, data, subset, na.action, ...)
```

lists arguments the function expects. If there is an equals sign, it has a default value that will be used if it is not specified.

You have to give it at least one vector of values ("x")

If you give it only x and nothing else, you get a one-sample test of the null hypothesis that the population mean of x (mu) is 0 and a two-sided alternative

Details

`alternative = "greater"` is the alternative that `x` has a larger mean than `y`. For the one-sample case: that the mean is positive.

If `paired` is `TRUE` then both `x` and `y` must be specified and they must be the same length. Missing values are silently removed (in pairs if `paired` is `TRUE`). If `var.equal` is `TRUE` then the pooled estimate of the variance is used. By default, if `var.equal` is `FALSE` then the variance is estimated separately for both groups and the Welch modification to the degrees of freedom is used.

If the input data are effectively constant (compared to the larger of the two means) an error is generated.

Value

A list with class "htest" containing the following components:

<code>statistic</code>	the value of the t-statistic.
<code>parameter</code>	the degrees of freedom for the t-statistic.
<code>p.value</code>	the p-value for the test.

```
Files
R: Stud
t.test
One Sample t-test
data: group1
t = 23.172, df = 9, p-value = 2.47e-09
alternative hypothesis: true mean is not equal to 0
95 percent confidence interval:
 45.38939 55.21061
sample estimates:
mean of x
 50.3
> t.test(group1, group2)
Welch Two Sample t-test
data: group1 and group2
t = -3.5069, df = 15.894, p-value = 0.002946
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 -21.665073 -5.334927
sample estimates:
mean of x mean of y
 50.3      63.8
> t.test()
Error in t.test.default() : argument "x" is missing, with no default
```

Details

`alternative = "greater"` is the alternative that `x` has a larger mean than `y`. For the one-sample case: that the mean is positive.

If `paired` is `TRUE` then both `x` and `y` must be specified and they must be the same length. Missing values are silently removed (in pairs if `paired` is `TRUE`). If `var.equal` is `TRUE` then the pooled estimate of the variance is used. By default, if `var.equal` is `FALSE` then the variance is estimated separately for both groups and the Welch modification to the degrees of freedom is used.

If the input data are effectively constant (compared to the larger of the two means) an error is generated.

Value

A list with class `"htest"` containing the following components:

<code>statistic</code>	the value of the t-statistic.
<code>parameter</code>	the degrees of freedom for the t-statistic.
<code>p.value</code>	the p-value for the test.

?t.test at prompt

Usage

```
t.test(x, ...)
```

```
## Default S3 method:
```

```
t.test(x, y = NULL,
       alternative = c("two.sided", "less", "greater"),
       mu = 0, paired = FALSE, var.equal = FALSE,
       conf.level = 0.95, ...)
```

```
## S3 method for class 'formula'
```

```
t.test(formula, data, subset, na.action, ...)
```

lists arguments the function expects. If there is an equals sign, it has a default value that will be used if it is not specified.

Default for `t.test` is for `var.equal` to be `FALSE` i.e. unequal variances, so if you give it 2 vectors of values and nothing else, it gives you a t-test with a correction factor for estimating separate variances in each group

```
> t.test(group1, group2, var.equal = TRUE)
```

Two Sample t-test

```
data: group1 and group2
t = -3.5069, df = 18, p-value = 0.002518
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 -21.587565 -5.412435
sample estimates:
mean of x mean of y
 50.3      63.8
```

```
> t.test(group1, group2)
```

Welch Two Sample t-test

```
data: group1 and group2
t = -3.5069, df = 15.894, p-value = 0.002946
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 -21.665073 -5.334927
sample estimates:
mean of x mean of y
 50.3      63.8
```

```
> t.test()
Error in t.test.default() : argument "x" is missing, with no default
```

?t.test at prompt

Usage

```
t.test(x, ...)
```

```
## Default S3 method:
```

```
t.test(x, y = NULL,
       alternative = c("two.sided", "less", "greater"),
       mu = 0, paired = FALSE, var.equal = FALSE,
       conf.level = 0.95, ...)
```

```
## S3 method for class 'formula'
```

```
t.test(formula, data, subset, na.action, ...)
```

lists arguments the function expects. If there is an equals sign, it has a default value that will be used if it is not specified.

If you set `var.equal` to `TRUE` then you get the more typical behavior of a t-test based on a pooled variance estimate (assume population variances of the two groups are the same)

`alternative = "greater"` is the alternative that `x` has a larger mean than `y`. For the one-sample case: that the mean is positive.

If `paired` is `TRUE` then both `x` and `y` must be specified and they must be the same length. Missing values are silently removed (in pairs if `paired` is `TRUE`). If `var.equal` is `TRUE` then the pooled estimate of the variance is used. By default, if `var.equal` is `FALSE` then the variance is estimated separately for both groups and the Welch modification to the degrees of freedom is used.

If the input data are effectively constant (compared to the larger of the two means) an error is generated.

Value

A list with class `"htest"` containing the following components:

<code>statistic</code>	the value of the t-statistic.
<code>parameter</code>	the degrees of freedom for the t-statistic.
<code>p.value</code>	the p-value for the test.


```
> my_data <- data.frame(group = c(rep("group1",10), rep("group2", 10)), values = c(group1, group2))
```

```
> my_data
```

```
  group values
```

```
1 group1    56
```

```
2 group1    46
```

```
3 group1    45
```

```
4 group1    42
```

```
5 group1    60
```

```
6 group1    45
```

```
7 group1    52
```

```
8 group1    59
```

```
9 group1    43
```

```
10 group1   55
```

```
11 group2    85
```

```
12 group2    61
```

```
13 group2    57
```

```
14 group2    53
```

```
15 group2    64
```

```
16 group2    58
```

```
17 group2    67
```

```
18 group2    54
```

```
19 group2    76
```

```
20 group2    63
```

```
> t.test(values ~ group, data = my_data, var.equal = TRUE)
```

Two Sample t-test

data: values by group

t = -3.5069, df = 18, p-value = 0.002518

alternative hypothesis: true difference in means between group group1 and group group2 is not equal to 0

95 percent confidence interval:

-21.587565 -5.412435

sample estimates:

mean in group group1 mean in group group2

50.3

63.8

parameter the degrees of freedom for the t-statistic.

p.value the p-value for the test.

prompt

Usage

```
t.test(x, ...)
```

```
## Default S3 method:
```

```
t.test(x, y = NULL,
       alternative = c("two.sided", "less", "greater"),
       mu = 0, paired = FALSE, var.equal = FALSE,
       conf.level = 0.95, ...)
```

```
## S3 method for class 'formula'
```

```
t.test(formula, data, subset, na.action, ...)
```

formula syntax: $y \sim x$

"y explained by x"

if you have data with group in one column and the value in the other, you can use `t.test` with "formula" syntax – gives same answer

Questions and discussion on "pre-test"

- confidence interval more informative than binary "significant" / "not significant" conclusion based on arbitrary cutoff
- interpretation can be nuanced – inference is based on a population parameter that is a fixed quantity, so probability is on the behavior of the confidence interval
- in practice, thinking of confidence interval as putting a probability on the population value is usually not too far off (but is wrong)



Questions and discussion on "pre-test"

```
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 -21.587565  -5.412435
sample estimates:
mean of x mean of y
   50.3    63.8
```

- in this case, the confidence interval is on the mean of group1 (x) minus the mean of group2 (y)
- the 95% CI doesn't contain 0, so we conclude the means of the 2 groups are "significantly" different at the .05 level
- the difference could be between 5.4 and 21.5



Permutation tests in R

