

# Introduction to R and Programming

Nathaniel E. Helwig

Assistant Professor of Psychology and Statistics  
University of Minnesota (Twin Cities)



Updated 04-Jan-2017

Copyright © 2017 by Nathaniel E. Helwig

# Outline of Notes

## 1) Introduction to R:

- Downloading R
- Basic calculations
- Using R functions
- Object classes in R

## 3) Basic Programming:

- Logical operators
- If/Else statements
- For loops
- While statements

## 2) Statistical Distributions in R:

- Overview
- Normal distribution
- Student's  $t$  distribution
- Common distributions

# Introduction to R

# R = Free and Open-Source Statistics

R is a **free** and **open-source** software environment for statistics.

- Created by Ross Ihaka and Robert Gentleman (at the University of Auckland, New Zealand)
- Based on S language created by John Chambers (at Bell Labs)
- Currently managed by The R Project for Statistical Computing  
<http://www.r-project.org/>

You can freely download R for various operating systems:

- Mac
- Windows
- Linux

# RStudio IDE

RStudio IDE is a **free** and **open-source** integrated development environment (IDE) for R.

- Basic R interface is a bit rough (particularly on Windows)
- RStudio offers a nice environment through which you can use R
- Freely available at <http://www.rstudio.com/>

You can freely download RStudio IDE for various operating systems:

- Mac
- Windows
- Linux

# R Console as a Calculator

## Addition and Subtraction

```
> 3+2  
[1] 5
```

```
> 3-2  
[1] 1
```

## Multiplication and Division

```
> 3*2  
[1] 6
```

```
> 3/2  
[1] 1.5
```

## Exponents in R

```
> 3^2  
[1] 9
```

```
> 2^3  
[1] 8
```

## Constants in R

```
> pi  
[1] 3.141593
```

```
> exp(1)  
[1] 2.718282
```

# Some Special Values in R

## Infinite Values

```
> Inf  
[1] Inf
```

```
> 1+Inf  
[1] Inf
```

## Machine Epsilon

```
> .Machine$double.eps  
[1] 2.220446e-16
```

```
> 0>.Machine$double.eps  
[1] FALSE
```

## Empty Values

```
> NULL  
NULL
```

```
> 1+NULL  
numeric(0)
```

## Missing Values

```
> NA  
[1] NA
```

```
> 1+NA  
[1] NA
```



# Storing and Manipulating Values in R

Define objects `x` and `y` with values of 3 and 2, respectively:

```
> x=3
```

```
> y=2
```

Some calculations with the defined objects `x` and `y`:

```
> x+y
```

```
[1] 5
```

```
> x*y
```

```
[1] 6
```

Warning: R is case sensitive, so `x` and `X` are not the same object.

# Function-Based Languages

R is a function-based language, where a “function” takes in some input  $\mathfrak{X}_I$  and creates some output  $\mathfrak{X}_O$ .

Vegas rules: what happens in a function, stays in a function

- Function only knows the input  $\mathfrak{X}_I$
- Function only creates the output  $\mathfrak{X}_O$

Each R function has a (unique) name, and the general syntax is

$$\mathfrak{X}_O = \text{fname}(\mathfrak{X}_I, \dots)$$

where `fname` is the function name, and  $\dots$  denotes additional inputs.

# Some Basic R Functions

## Combine

```
> c(1, 3, -2)
[1] 1 3 -2
```

```
> c("a", "a", "b", "b", "a")
[1] "a" "a" "b" "b" "a"
```

## Sum and Mean

```
> sum(c(1, 3, -2))
[1] 2
```

```
> mean(c(1, 3, -2))
[1] 0.6666667
```

## Variance and Std. Dev.

```
> var(c(1, 3, -2))
[1] 6.333333
```

```
> sd(c(1, 3, -2))
[1] 2.516611
```

## Minimum and Maximum

```
> min(c(1, 3, -2))
[1] -2
```

```
> max(c(1, 3, -2))
[1] 3
```

# Some More R Functions

Define objects `x` and `y`:

```
> x=c(1,3,4,6,8)
> y=c(2,3,5,7,9)
```

Calculate the correlation:

```
> cor(x,y)
[1] 0.988765
```

Calculate the covariance:

```
> cov(x,y)
[1] 7.65
```

Combine as columns

```
> cbind(x,y)
      x y
[1,] 1 2
[2,] 3 3
[3,] 4 5
[4,] 6 7
[5,] 8 9
```

Combine as rows

```
> rbind(x,y)
      [,1] [,2] [,3] [,4] [,5]
x        1    3    4    6    8
y        2    3    5    7    9
```

# Object-Oriented Style Programming

R is an object-oriented language, where an “object” is a general term.

Any R object  $x$  has an associated “class”, which indicates the type of object that  $x$  represents.

Some R functions are only defined for a particular class of input  $x$ .

Other R functions perform different operations depending on the class of the input object  $x$ .

# Some Basic R Classes

numeric **class**:

```
> x=c(1,3,-2)
> x
[1] 1 3 -2
> class(x)
[1] "numeric"
```

integer **class**:

```
> x=c(1L,3L,-2L)
> x
[1] 1 3 -2
> class(x)
[1] "integer"
```

character **class**:

```
> x=c("a","a","b")
> x
[1] "a" "a" "b"
> class(x)
[1] "character"
```

factor **class**:

```
> x=factor(c("a","a","b"))
> x
[1] a a b
Levels: a b
> class(x)
[1] "factor"
```

# Some More R Classes

`matrix` **class**:

```
> x=c(1,3,-2)
> y=c(2,0,7)
> z=cbind(x,y)
> z
```

```
      x y
[1,]  1 2
[2,]  3 0
[3,] -2 7
```

```
> class(z)
[1] "matrix"
```

`data.frame` **class**:

```
> x=c(1,3,-2)
> y=c("a","a","b")
> z=data.frame(x,y)
> z
```

```
      x y
1     1 a
2     3 a
3    -2 b
```

```
> class(z)
[1] "data.frame"
```

# Class-Customized R Functions

Many functions in R are “class-customized”, i.e., they execute different code depending the on class of the input object  $\mathcal{X}$ .

One simple example (that we’ve already seen) is the `print` function:

```
> x=c(1,3,-2)
> y=factor(c("a","a","b"))
> print(x)
[1] 1 3 -2
> print(y)
[1] a a b
Levels: a b
```



# Class-Customized R Functions (continued)

Another simple example is the `summary` function:

```
> x=c(1,3,-2)
> y=factor(c("a","a","b"))
> summary(x)
  Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
-2.0000 -0.5000   1.0000   0.6667   2.0000   3.0000
> summary(y)
a b
2 1
```

# Class-Customized R Functions (continued)

Some R functions only work on particular object classes (e.g., `range`):

```
> x=c(1,3,-2)
```

```
> y=factor(c("a","a","b"))
```

```
> range(x)
```

```
[1] -2  3
```

```
> range(y)
```

```
Error in Summary.factor(c(1L, 1L, 2L), na.rm = FALSE)
  range not meaningful for factors
```

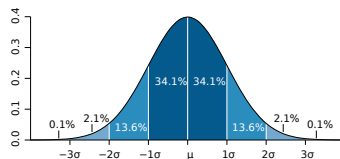
# Statistical Distributions in R

# Statistical Distributions: Summary

When working with different statistical distributions, we often want to make probabilistic statements based on the distribution.

We typically want to know one of three things:

- The density (pdf) value at a particular value of  $x$
- The distribution (cdf) value at a particular value of  $x$
- The quantile ( $x$ ) value corresponding to a particular probability



Statistical tables used to be printed in book appendices:

Entry is area  $A$  under the standard normal curve from  $-\infty$  to  $z(A)$



x	00	01	02	03	04	05	06	07	08	09
0	.5000	.5040	.5080	.5120	.5160	.5199	.5239	.5279	.5319	.5359
1	.5398	.5438	.5478	.5517	.5557	.5596	.5636	.5675	.5714	.5753
2	.5793	.5832	.5871	.5910	.5949	.5987	.6026	.6065	.6104	.6143
3	.6182	.6221	.6260	.6298	.6337	.6376	.6415	.6454	.6493	.6532
4	.6571	.6610	.6648	.6687	.6726	.6765	.6804	.6843	.6882	.6921
5	.6959	.6998	.7037	.7076	.7115	.7154	.7193	.7232	.7271	.7310
6	.7349	.7388	.7427	.7466	.7505	.7544	.7583	.7622	.7661	.7700
7	.7739	.7778	.7817	.7856	.7895	.7934	.7973	.8012	.8051	.8090
8	.8129	.8168	.8207	.8246	.8285	.8324	.8363	.8402	.8441	.8480
9	.8519	.8558	.8597	.8636	.8675	.8714	.8753	.8792	.8831	.8870
10	.8909	.8948	.8987	.9026	.9065	.9104	.9143	.9182	.9221	.9260
11	.9299	.9338	.9377	.9416	.9455	.9494	.9533	.9572	.9611	.9650
12	.9689	.9728	.9767	.9806	.9845	.9884	.9923	.9962	.9999	
13	.0000	.0039	.0078	.0117	.0156	.0195	.0234	.0273	.0312	.0351
14	.0390	.0429	.0468	.0507	.0546	.0585	.0624	.0663	.0702	.0741
15	.0780	.0819	.0858	.0897	.0936	.0975	.1014	.1053	.1092	.1131
16	.1170	.1209	.1248	.1287	.1326	.1365	.1404	.1443	.1482	.1521
17	.1560	.1599	.1638	.1677	.1716	.1755	.1794	.1833	.1872	.1911
18	.1950	.1989	.2028	.2067	.2106	.2145	.2184	.2223	.2262	.2301
19	.2340	.2379	.2418	.2457	.2496	.2535	.2574	.2613	.2652	.2691
20	.2730	.2769	.2808	.2847	.2886	.2925	.2964	.3003	.3042	.3081
21	.3120	.3159	.3198	.3237	.3276	.3315	.3354	.3393	.3432	.3471
22	.3510	.3549	.3588	.3627	.3666	.3705	.3744	.3783	.3822	.3861
23	.3900	.3939	.3978	.4017	.4056	.4095	.4134	.4173	.4212	.4251
24	.4290	.4329	.4368	.4407	.4446	.4485	.4524	.4563	.4602	.4641
25	.4680	.4719	.4758	.4797	.4836	.4875	.4914	.4953	.4992	.5031
26	.5070	.5109	.5148	.5187	.5226	.5265	.5304	.5343	.5382	.5421
27	.5460	.5499	.5538	.5577	.5616	.5655	.5694	.5733	.5772	.5811
28	.5850	.5889	.5928	.5967	.6006	.6045	.6084	.6123	.6162	.6201
29	.6240	.6279	.6318	.6357	.6396	.6435	.6474	.6513	.6552	.6591
30	.6630	.6669	.6708	.6747	.6786	.6825	.6864	.6903	.6942	.6981
31	.7020	.7059	.7098	.7137	.7176	.7215	.7254	.7293	.7332	.7371
32	.7410	.7449	.7488	.7527	.7566	.7605	.7644	.7683	.7722	.7761
33	.7800	.7839	.7878	.7917	.7956	.7995	.8034	.8073	.8112	.8151
34	.8190	.8229	.8268	.8307	.8346	.8385	.8424	.8463	.8502	.8541
35	.8580	.8619	.8658	.8697	.8736	.8775	.8814	.8853	.8892	.8931
36	.8970	.9009	.9048	.9087	.9126	.9165	.9204	.9243	.9282	.9321
37	.9360	.9399	.9438	.9477	.9516	.9555	.9594	.9633	.9672	.9711
38	.9750	.9789	.9828	.9867	.9906	.9945	.9984			

Cumulative probability A:	.90	.95	.975	.98	.99	.995	.999
z(A):	1.282	1.645	1.960	2.054	2.326	2.576	3.090

# Statistical Distributions: R Functions

R has functions for obtaining density, distribution, and quantile values.

The general naming structure of the relevant R functions is...

- `dname` calculates density (pdf) value at input quantile
- `pname` calculates distribution (cdf) value at input quantile
- `qname` calculates quantile value at input probability
- `rname` generates random sample of input size

Note that `name` represents the name of the given distribution.

# Normal Distribution: Overview

The relevant functions for the normal distribution are...

- `dnorm` calculates density (pdf) value at input quantile
- `pnorm` calculates distribution (cdf) value at input quantile
- `qnorm` calculates quantile value at input probability
- `rnorm` generates random sample of input size

In addition to the input quantile (or probability or size) value, you can input the `mean` and `sd` (standard deviation) of the variable.

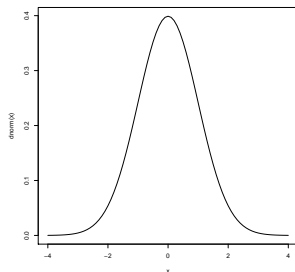
# Normal: Density Function

Standard normal density:

```
> dnorm(-4)
[1] 0.0001338302
> dnorm(-2)
[1] 0.05399097
> dnorm(0)
[1] 0.3989423
> dnorm(2)
[1] 0.05399097
> dnorm(4)
[1] 0.0001338302
```

Plot standard normal density:

```
> x=seq(-4, 4, by=.1)
> plot(x, dnorm(x), type="l")
```





# Normal: Density Function (continued)

Normal density with different mean and variance ( $\mu = 1$  and  $\sigma^2 = 2$ ):

```
> dnorm(-3, mean=1, sd=sqrt(2))
```

```
[1] 0.005166746
```

```
> dnorm(-1, mean=1, sd=sqrt(2))
```

```
[1] 0.1037769
```

```
> dnorm(1, mean=1, sd=sqrt(2))
```

```
[1] 0.2820948
```

```
> dnorm(3, mean=1, sd=sqrt(2))
```

```
[1] 0.1037769
```

```
> dnorm(5, mean=1, sd=sqrt(2))
```

```
[1] 0.005166746
```

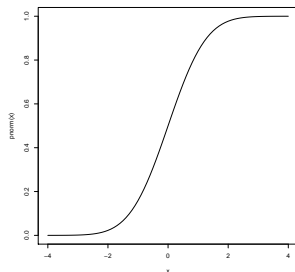
# Normal: Distribution Function

Standard normal cdf:

```
> pnorm(-4)
[1] 3.167124e-05
> pnorm(-2)
[1] 0.02275013
> pnorm(0)
[1] 0.5
> pnorm(2)
[1] 0.9772499
> pnorm(4)
[1] 0.9999683
```

Plot standard normal cdf:

```
> x=seq(-4, 4, by=.1)
> plot(x, pnorm(x), type="l")
```



# Normal: Distribution Function (continued)

Normal cdf with different mean and variance ( $\mu = 1$  and  $\sigma^2 = 2$ ):

```
> pnorm(-3, mean=1, sd=sqrt(2))  
[1] 0.002338867  
> pnorm(-1, mean=1, sd=sqrt(2))  
[1] 0.0786496  
> pnorm(1, mean=1, sd=sqrt(2))  
[1] 0.5  
> pnorm(3, mean=1, sd=sqrt(2))  
[1] 0.9213504  
> pnorm(5, mean=1, sd=sqrt(2))  
[1] 0.9976611
```

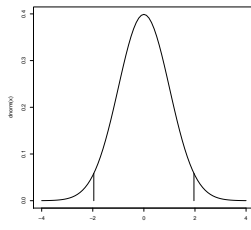
# Normal: Quantile Function

Standard normal quantiles:

```
> qnorm(.005)
[1] -2.575829
> qnorm(.025)
[1] -1.959964
> qnorm(.5)
[1] 0
> qnorm(.975)
[1] 1.959964
> qnorm(.995)
[1] 2.575829
```

Plot standard normal quantiles:

```
x=seq(-4, 4, by=.1)
plot(x, dnorm(x), type="l")
qx=qnorm(.025)
lines(x=rep(qx, 2),
      y=c(0, dnorm(qx)))
lines(x=rep(-qx, 2),
      y=c(0, dnorm(-qx)))
```



# Normal: Quantile Function (continued)

Normal quantiles with different mean and variance ( $\mu = 1$  and  $\sigma^2 = 2$ ):

```
> qnorm(.005, mean=1, sd=sqrt(2))
```

```
[1] -2.642773
```

```
> qnorm(.025, mean=1, sd=sqrt(2))
```

```
[1] -1.771808
```

```
> qnorm(.5, mean=1, sd=sqrt(2))
```

```
[1] 1
```

```
> qnorm(.975, mean=1, sd=sqrt(2))
```

```
[1] 3.771808
```

```
> qnorm(.995, mean=1, sd=sqrt(2))
```

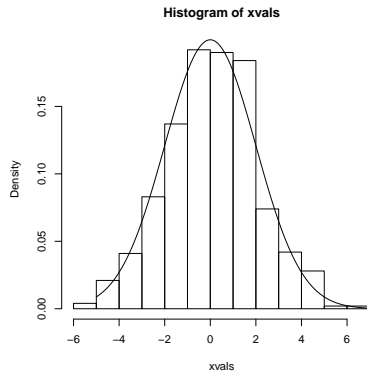
```
[1] 4.642773
```

# Simulating Normal Data in R

For each distribution, use `rname` to simulate from `name` distribution.

For example, to simulate normal

```
> set.seed(12345)
> xvals=rnorm(1000,mean=0,sd=2)
> xseq=seq(-5,7,l=100)
> hist(xvals,freq=FALSE)
> lines(xseq,dnorm(xseq,sd=2))
```



# Testing Normality in R

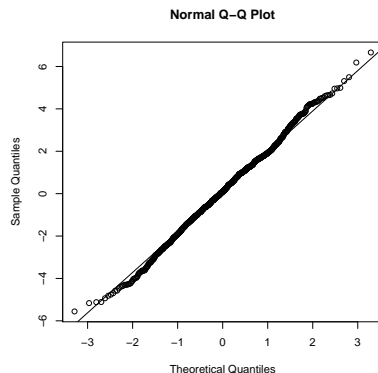
Use `qqnorm` and `qqline` to make Q-Q plot and `shapiro.test` to perform Shapiro-Wilk normality test.

For example, to test normality

```
> set.seed(12345)
> xvals=rnorm(1000,mean=0,sd=2)
> qqnorm(xvals)
> qqline(xvals)
> shapiro.test(xvals)
```

Shapiro-Wilk normality test

```
data:  xvals
W = 0.9978, p-value = 0.1988
```



# Student's $t$ Distribution: Overview

Family of real-valued continuous distributions that depends on the parameter  $\nu > 0$ , which is the **degrees of freedom**.

We encounter  $t$  distribution when estimating  $\mu$  (the mean of a normal variable) with  $\sigma^2$  (the variance of the normal variable) unknown.

Called “Student’s”  $t$  because of William Gosset. . .

- Worked for Guinness Brewery (Dublin, Ireland) in early 1900s
- Published paper under pseudonym “Student” because Guinness did not allow employees to publish scientific papers

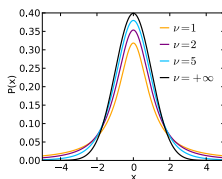


# Student's $t$ Distribution: Properties

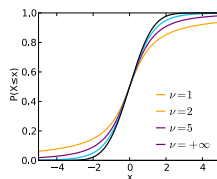
Like the standard normal distribution, the  $t$  distribution is bell-shaped and symmetric around zero.

For small  $\nu$ , the  $t$  distribution has heavy tails; as  $\nu \rightarrow \infty$ ,  $t$  distribution approaches standard normal distribution.

Helpful figures of  $t$  distribution pdfs and cdfs:



[http://en.wikipedia.org/wiki/File:Student\\_t\\_pdf.svg](http://en.wikipedia.org/wiki/File:Student_t_pdf.svg)



[http://en.wikipedia.org/wiki/File:Student\\_t\\_cdf.svg](http://en.wikipedia.org/wiki/File:Student_t_cdf.svg)

# Student's $t$ Distribution: R Functions

The relevant functions for the  $t$  distribution are...

- `dt` calculates density (pdf) value at input quantile
- `pt` calculates distribution (cdf) value at input quantile
- `qt` calculates quantile value at input probability
- `rt` generates random sample of input size

In addition to the input quantile (or probability or size) value, you can input the `df` (degrees of freedom) and `ncp` (non-centrality parameter).

- We will not discuss non-central  $t$  distributions
- You only need to worry about the `df` input

# Student's $t$ Distribution: Example Code

Student's  $t$  pdf (at  $x = 0$ ):

```
> dt(0, df=1)
[1] 0.3183099
> dt(0, df=10)
[1] 0.3891084
> dt(0, df=100)
[1] 0.3979462
```

Student's  $t$  cdf (at  $x = 0$ ):

```
> pt(0, df=1)
[1] 0.5
> pt(0, df=10)
[1] 0.5
> pt(0, df=100)
[1] 0.5
```

Student's  $t$  quantiles (at  $p = .975$ ):

```
> qt(.975, df=1)
[1] 12.7062
> qt(.975, df=10)
[1] 2.228139
> qt(.975, df=100)
[1] 1.983972
```

Student's  $t$  quantiles (at  $p = .995$ ):

```
> qt(.995, df=1)
[1] 63.65674
> qt(.995, df=10)
[1] 3.169273
> qt(.995, df=100)
[1] 2.625891
```

# One Sample $t$ Test: Overview

Suppose  $x_i \stackrel{\text{iid}}{\sim} N(\mu, \sigma^2)$  and want to test  $H_0 : \mu = \mu_0$  versus  $H_1 : \mu \neq \mu_0$

Assuming  $\sigma$  is unknown, use the one-sample Student's  $t$  test statistic:

$$T = \frac{\bar{x} - \mu_0}{s/\sqrt{n}} \sim t_{n-1}$$

where  $\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$  and  $s = \sqrt{\frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}}$

100(1 -  $\alpha$ )% confidence interval (CI) for  $\mu$  is given by

$$\bar{x} \pm t_{n-1}^{(\alpha/2)}(s/\sqrt{n})$$

where  $t_{n-1}^{(\alpha/2)}$  is critical  $t_{n-1}$  value such that  $P(T > t_{n-1}^{(\alpha/2)}) = \alpha/2$ .

# One Sample $t$ Test: Example

A store sells “16-ounce” boxes of Captain Crisp cereal. A random sample of 9 boxes was taken and weighed. The results were

15.5 16.2 16.1 15.8 15.6 16.0 15.8 15.9 16.2

ounces. Assume the weight of cereal in a box is normally distributed.

The sample mean and variance are given by

$$\bar{x} = (1/n) \sum_{i=1}^n x_i = (1/9)(15.5 + \cdots + 16.2) = (1/9)(143.1) = 15.9$$

$$\begin{aligned} s^2 &= (n-1)^{-1} \sum_{i=1}^n (x_i - \bar{x})^2 = (n-1)^{-1} \left[ \sum_{i=1}^n x_i^2 - n\bar{x}^2 \right] \\ &= (1/8) \left[ 2275.79 - 9(15.9^2) \right] = (1/8)(0.5) = 0.0625 \end{aligned}$$

## One Sample $t$ Test: Example (continued)

$t_8^{(.025)} = 2.306$ , so the 95% CI for the average weight of a cereal box is:  
 $15.9 \pm 2.306\sqrt{0.0625/9} = [15.708; 16.092]$

The company that makes Captain Crisp cereal claims that the average weight of its box is at least 16 ounces. Use a 0.05 level of significance to test the company's claim. What is the p-value of this test?

To test  $H_0 : \mu \geq 16$  versus  $H_1 : \mu < 16$ , the test statistic is

$$T = \frac{15.9 - 16}{\sqrt{0.0625/9}} = -1.2$$

We know that  $T \sim t_8$ , so we have that  $P(T < -1.2) = 0.1322336$ .  
Therefore, we retain  $H_0$  at the  $\alpha = .05$  level.

# One Sample $t$ Test: R Code

```
> x=c(15.5, 16.2, 16.1, 15.8, 15.6, 16.0, 15.8, 15.9, 16.2)
> mean(x)
[1] 15.9
> sd(x)
[1] 0.25
> var(x)
[1] 0.0625
> t.test(x)
```

One Sample t-test

```
data:  x
t = 190.8, df = 8, p-value = 6.372e-16
alternative hypothesis: true mean is not equal to 0
95 percent confidence interval:
 15.70783 16.09217
sample estimates:
mean of x
 15.9
```

# One Sample $t$ Test: R Code (continued)

```
> t.test(x,mu=16,alternative="less",conf.level=.95)
```

```
One Sample t-test
```

```
data:  x
t = -1.2, df = 8, p-value = 0.1322
alternative hypothesis: true mean is less than 16
95 percent confidence interval:
    -Inf 16.05496
sample estimates:
mean of x
    15.9
```



# Two Sample $t$ Test: Overview

Suppose  $x_i \stackrel{\text{iid}}{\sim} N(\mu_x, \sigma^2)$  and  $y_i \stackrel{\text{iid}}{\sim} N(\mu_y, \sigma^2)$

Want to test  $H_0 : \mu_x - \mu_y = \mu_0$  versus  $H_1 : \mu_x - \mu_y \neq \mu_0$

Assuming  $\sigma$  is unknown, use the two-sample Student's  $t$  test statistic:

$$T = \frac{(\bar{x} - \bar{y}) - \mu_0}{s_p \sqrt{\frac{1}{n} + \frac{1}{m}}} \sim t_{n+m-2}$$

where  $\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$ ,  $\bar{y} = \frac{\sum_{i=1}^m y_i}{m}$ , and  $s_p^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2 + \sum_{i=1}^m (y_i - \bar{y})^2}{n+m-2}$

100(1 -  $\alpha$ )% confidence interval (CI) for  $\mu_x - \mu_y$  is given by

$$(\bar{x} - \bar{y}) \pm t_{n+m-2}^{(\alpha/2)} \left( s_p \sqrt{\frac{1}{n} + \frac{1}{m}} \right)$$

where  $t_{n+m-2}^{(\alpha/2)}$  is critical  $t_{n+m-2}$  value such that  $P(T > t_{n+m-2}^{(\alpha/2)}) = \alpha/2$ .

## Two Sample $t$ Test: Example

Assume that the distributions of  $X$  and  $Y$  are  $N(\mu_1, \sigma^2)$  and  $N(\mu_2, \sigma^2)$ , respectively. Given the  $n = 6$  observations of  $X$ ,

70,      82,      78,      74,      94,      82

and the  $m = 8$  observations of  $Y$ ,

64,      72,      60,      76,      72,      80,      84,      68

find the p-value for the test  $H_0 : \mu_1 = \mu_2$  versus  $H_1 : \mu_1 > \mu_2$ .

## Two Sample $t$ Test: Example (continued)

First, note that the sample means and variances are given by

$$\bar{x} = (1/6) \sum_{i=1}^6 x_i = (1/6)480 = 80$$

$$\bar{y} = (1/8) \sum_{i=1}^8 y_i = (1/8)576 = 72$$

$$s_x^2 = (1/5) \sum_{i=1}^6 (x_i - \bar{x})^2 = (1/5)344 = 68.8$$

$$s_y^2 = (1/7) \sum_{i=1}^8 (y_i - \bar{y})^2 = (1/7)448 = 64$$

which implies that the pooled variance estimate is given by

$$\begin{aligned} s_p^2 &= \frac{(n-1)s_x^2 + (m-1)s_y^2}{n+m-2} \\ &= \frac{344 + 448}{12} \\ &= 66 \end{aligned}$$

## Two Sample $t$ Test: Example (continued)

Thus, the relevant  $t$  test statistic is given by

$$\begin{aligned} T &= \frac{(\bar{X} - \bar{Y}) - \mu_0}{s_p \sqrt{\frac{1}{n} + \frac{1}{m}}} \\ &= \frac{(80 - 72) - 0}{\sqrt{66} \sqrt{\frac{1}{6} + \frac{1}{8}}} \\ &= 1.82337 \end{aligned}$$

Note that  $T \sim t_{12}$ , so the corresponding p-value is

$$P(T > 1.82337) = 0.04661955$$

Therefore, we reject  $H_0$  at the  $\alpha = .05$  level.

# Two Sample $t$ Test: R Code

```
> x=c(70, 82, 78, 74, 94, 82)
> y=c(64, 72, 60, 76, 72, 80, 84, 68)
> t.test(x,y,alternative="greater",var.equal=TRUE)
```

Two Sample  $t$ -test

```
data:  x and y
t = 1.8234, df = 12, p-value = 0.04662
alternative hypothesis: true difference in means is greater than 0
95 percent confidence interval:
 0.1802451      Inf
sample estimates:
mean of x mean of y
    80      72
```

# Chi-Squared Distribution: Overview

Family of positive real-valued continuous distributions that depends on the parameter  $k > 0$ , which is the **degrees of freedom**.

If  $Z_1, \dots, Z_k$  are iid  $N(0, 1)$ , then  $Q = (\sum_{i=1}^k Z_i^2) \sim \chi_k^2$

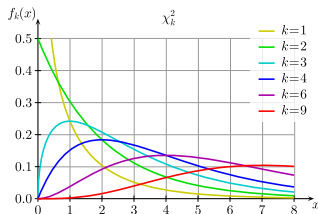
- iid: independent identically distributed
- $\chi_k^2$  denotes a chi-squared distribution with  $k$  degrees of freedom

# Chi-Squared Distribution: Properties

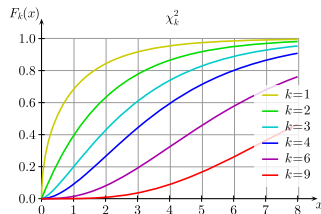
Chi-squared variable must be nonnegative (squared normal variable).

$\chi_k^2$  distribution takes a variety of different shapes depending on  $k$ .

Helpful figures of  $\chi_k^2$  distribution pdfs and cdfs:



[http://en.wikipedia.org/wiki/File:Chi-square\\_pdf.svg](http://en.wikipedia.org/wiki/File:Chi-square_pdf.svg)



[http://en.wikipedia.org/wiki/File:Chi-square\\_cdf.svg](http://en.wikipedia.org/wiki/File:Chi-square_cdf.svg)

# Chi-Squared Distribution: R Functions

The relevant functions for the  $\chi_k^2$  distribution are...

- `dchisq` calculates density (pdf) value at input quantile
- `pchisq` calculates distribution (cdf) value at input quantile
- `qchisq` calculates quantile value at input probability
- `rchisq` generates random sample of input size

In addition to the input quantile (or probability or size) value, you can input the `df` (degrees of freedom) and `ncp` (non-centrality parameter).

- We will not discuss non-central  $\chi_k^2$  distributions
- You only need to worry about the `df` input



# Chi-Squared Distribution: Example Code

$\chi_k^2$  pdf (at  $x = 1$ ):

```
> dchisq(1,df=1)
[1] 0.2419707
> dchisq(1,df=10)
[1] 0.0007897535
> dchisq(1,df=100)
[1] 8.856214e-79
```

$\chi_k^2$  cdf (at  $x = 1$ ):

```
> pchisq(1,df=1)
[1] 0.6826895
> pchisq(1,df=10)
[1] 0.0001721156
> pchisq(1,df=100)
[1] 1.788777e-80
```

$\chi_k^2$  quantiles (at  $p = .975$ ):

```
> qchisq(.975,df=1)
[1] 5.023886
> qchisq(.975,df=10)
[1] 20.48318
> qchisq(.975,df=100)
[1] 129.5612
```

$\chi_k^2$  quantiles (at  $p = .995$ ):

```
> qchisq(.995,df=1)
[1] 7.879439
> qchisq(.995,df=10)
[1] 25.18818
> qchisq(.995,df=100)
[1] 140.1695
```

# F Distribution: Overview

Family of positive real-valued continuous distributions that depends on the parameters  $k_1, k_2 > 0$ , which are the **numerator and denominator degrees of freedom**, respectively.

If  $Q_1 \sim \chi_{k_1}^2$  and  $Q_2 \sim \chi_{k_2}^2$  are independent, then  $F = \frac{Q_1/k_1}{Q_2/k_2} \sim F_{k_1, k_2}$

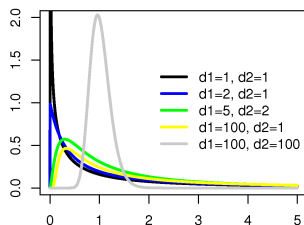
- independent:  $Q_1$  and  $Q_2$  are statistically independent
- $F_{k_1, k_2}$  denotes an  $F$  distribution with  $k_1$  numerator degrees of freedom and  $k_2$  denominator degrees of freedom

# F Distribution: Properties

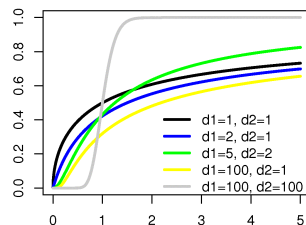
$F$  variables must be nonnegative (ratio of scaled chi-squared).

$F$  distribution takes a variety of different shapes depending on  $k_1, k_2$ .

Helpful figures of  $F$  distribution pdfs and cdfs:



[http://en.wikipedia.org/wiki/File:F\\_distributionPDF.png](http://en.wikipedia.org/wiki/File:F_distributionPDF.png)



[http://en.wikipedia.org/wiki/File:F\\_distributionCDF.png](http://en.wikipedia.org/wiki/File:F_distributionCDF.png)

# F Distribution: R Functions

The relevant functions for the  $F_{k_1, k_2}$  distribution are...

- `df` calculates density (pdf) value at input quantile
- `pf` calculates distribution (cdf) value at input quantile
- `qf` calculates quantile value at input probability
- `rf` generates random sample of input size

In addition to the input quantile/probability/size, you can input `df1` and `df2` (degrees of freedom), and `ncp` (non-centrality parameter).

- We will not discuss non-central  $F_{k_1, k_2}$  distributions
- You only need to worry about the `df1` and `df2` inputs

# F Distribution: Example Code

$F_{k_1, k_2}$  pdf (at  $x = 1$ ):

```
> df(1, df1=1, df2=1)
[1] 0.1591549
> df(1, df1=1, df2=10)
[1] 0.230362
> df(1, df1=10, df2=10)
[1] 0.6152344
```

$F_{k_1, k_2}$  cdf (at  $x = 1$ ):

```
> pf(1, df1=1, df2=1)
[1] 0.5
> pf(1, df1=1, df2=10)
[1] 0.6591069
> pf(1, df1=10, df2=10)
[1] 0.5
```

$F_{k_1, k_2}$  quantiles (at  $p = .975$ ):

```
> qf(.975, df1=1, df2=1)
[1] 647.789
> qf(.975, df1=1, df2=10)
[1] 6.936728
> qf(.975, df1=10, df2=10)
[1] 3.716792
```

$F_{k_1, k_2}$  quantiles (at  $p = .995$ ):

```
> qf(.995, df1=1, df2=1)
[1] 16210.72
> qf(.995, df1=1, df2=10)
[1] 12.82647
> qf(.995, df1=10, df2=10)
[1] 5.846678
```

# Other Distributions in R

R has many more distributions that we will not discuss, e.g.:

- Beta distribution (`dbeta`, `pbeta`, `qbeta`, `rbeta`)
- Binomial distribution (`dbinom`, `binom`, `qbinom`, `rbinom`)
- Exponential distribution (`dexp`, `pexp`, `qexp`, `rexp`)
- Gamma distribution (`dgamma`, `pgamma`, `qgamma`, `rgamma`)
- Log Normal distribution (`dlnorm`, `plnorm`, `qlnorm`, `rlnorm`)
- Negative Binomial (`dnbinom`, `pnbinom`, `qnbinom`, `rnbinom`)
- Poisson distribution (`dpois`, `ppois`, `qpois`, `rpois`)
- Uniform distribution (`dunif`, `punif`, `qunif`, `runif`)

Note: all R distributions follow the same naming convention.

# Basic Programming

# Logical Operators: Overview

Logical operators derive from Boolean algebra, where values of variables are either `TRUE` or `FALSE`.

We use logical operators to execute different code depending on whether a condition is met.

Logical operators are used within *many* R functions, so an understanding of logical operators is crucial to understanding R code.



# Logical Operators: R Syntax

Operator	Summary
<	Less than
>	Greater than
<=	Less than or equal to
>=	Greater than or equal to
==	Equal to
!=	Not equal to
!x	NOT x
x y	x OR y
x&y	x AND y

# Logical Operators: Example

Define objects  $x$  and  $y$ :

```
> x=y=10
```

Less than:

```
> x<y
```

```
[1] FALSE
```

```
> x<=y
```

```
[1] TRUE
```

Greater than:

```
> x>y
```

```
[1] FALSE
```

```
> x>=y
```

```
[1] TRUE
```

Equal to (not equal to):

```
> x==y
```

```
[1] TRUE
```

```
> x!=y
```

```
[1] FALSE
```

OR and AND:

```
> x=10
```

```
> y=11
```

```
> (x<11 | y<11)
```

```
[1] TRUE
```

```
> (x<11 & y<11)
```

```
[1] FALSE
```

# If/Else Statements: Overview

If/Else statements are fundamental in any programming language.

We use if/else statements (in combination with logical operators) to execute different code depending on whether a condition is met.

If/Else statements always appear with logical operators.

# If/Else Statements: R Syntax

General if/else syntax:

```
if (...) {  
  
    some R code  
  
} else {  
  
    more R code  
  
}
```

Nested if/else syntax:

```
if (...) {  
  
    some R code  
  
} else if (...) {  
  
    more R code  
  
} else {  
  
    even more R code  
  
}
```

# If/Else Statements: Example

```
> x=10
> if(x>5) {
+   x=x/2
+   y=2*x
+ } else {
+   x=x*2
+   y=x
+ }
> x
[1] 5
> y
[1] 10
```

```
> x=4
> if(x>5) {
+   x=x/2
+   y=2*x
+ } else {
+   x=x*2
+   y=x
+ }
> x
[1] 8
> y
[1] 8
```

Note: the + signs are NOT part of the R code; these are included by R when entering multiline statements.

# If/Else Statements: Example (continued)

To be more efficient, we could write an R function:

```
> myfun<-function(x) {  
+   if(x>5) {  
+     x=x/2  
+     y=2*x  
+   } else {  
+     x=x*2  
+     y=x  
+   }  
+   list(x=x,y=y)  
+ }  
  
> class(myfun)  
[1] "function"  
  
> myfun(10)  
$x  
[1] 5  
  
$y  
[1] 10  
  
> myfun(4)  
$x  
[1] 8  
  
$y  
[1] 8
```

# For Loops: Overview

For loops (or do loops) are fundamental in any programming language.

We use for loops to execute the same code repeatedly with the loop index changing at each step of the loop.

Warning: for loops in R can be slow; vectorize your code if possible!

# For Loops: Syntax

```
for(j in J){  
    some R code depending on j  
}
```

Note:  $j$  is the loop index and  $J$  is the index set.



# For Loops: Example

## For loop version:

```
> x=11:15
> x
[1] 11 12 13 14 15
> for(idx in 1:5){
+   x[idx]=x[idx]+1
+ }
> x
[1] 12 13 14 15 16
```

## Vectorized version:

```
> x=11:15
> x
[1] 11 12 13 14 15
> x=x+1
> x
[1] 12 13 14 15 16
```

# While Statements: Overview

While statements are fundamental in any programming language.

We use while statements (in combination with logical operators) to execute the same code repeatedly until some condition is met.

While statements always appear with logical operators.

# While Statements: Syntax

```
while(...) {  
  
    some R code  
  
}
```

Note: keeps repeating R code until logical statement . . . is FALSE

# While Statements: Example

Simple while statement:

```
> x=80
> iter=0
> while(x<100){
+   x=x+sqrt(x)/10
+   iter=iter+1
+ }
> x
[1] 100.8293
> iter
[1] 22
```

Another while statement:

```
> x=80
> iter=0
> while(x<100 & iter<20){
+   x=x+sqrt(x)/10
+   iter=iter+1
+ }
> x
[1] 98.83599
> iter
[1] 20
```

# While Statements: Example (continued)

Improper while statement:

```
> iter=0
> while(x<100){
+     x=x-sqrt(x)/10
+     iter=iter+1
+ }
```

```
Error in while (x < 100) {
  : missing value where TRUE/FALSE needed
In addition: Warning message:
In sqrt(x) : NaNs produced
```

**Note:** we get error message because  $x$  becomes negative, so we get NaN when we take the square-root.

# While Statements: Example (continued)

Infinite while statement:

```
> x=80
> iter=0
> while (x<100) {
+     x=x-x/10
+     iter=iter+1
+ }
```

Note: while statement will run infinitely (until we manually stop it) because logical statement is always true (i.e.,  $x < 100$  always).