

Statistics in R

James McDermott

NUI Galway



NUI Galway
OÉ Gaillimh

Programming and Tools for Artificial Intelligence

Dr James McDermott

Statistics in R

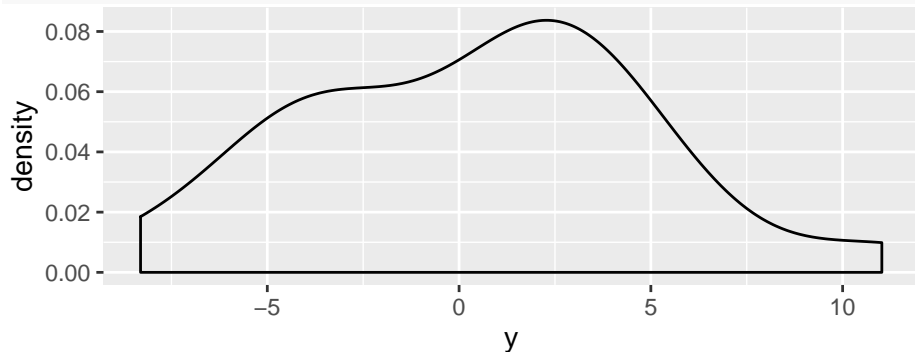
Load Tidyverse as usual

```
library(tidyverse)
```

```
## -- Attaching packages -----  
## v ggplot2 3.2.1      v purrr  0.3.2  
## v tibble  2.1.3      v dplyr  0.8.3  
## v tidyr   1.0.0      v stringr 1.4.0  
## v readr   1.3.1      v forcats 0.4.0  
## -- Conflicts -----  
## x dplyr::filter() masks stats::filter()  
## x dplyr::lag()    masks stats::lag()
```

Random numbers

```
x1 <- runif(20, min=0, max=2) # random uniform with bounds  
x2 <- rnorm(20) # random normal with mean 0, sd 1  
y <- x1 + x2 * rnorm(20, mean=5, sd=2)  
ggplot(tibble(y), aes(x=y)) + geom_density()
```



Basic statistics

```
for (f in c(min, max, mean, median, sd, var, IQR, mad)) {  
  print(f(y))  
}
```

```
## [1] -8.304924
```

```
## [1] 11.02246
```

```
## [1] 0.4374315
```

```
## [1] 1.325795
```

```
## [1] 4.580337
```

```
## [1] 20.97949
```

```
## [1] 5.971987
```

```
## [1] 4.859902
```

More data summaries

```
for (f in c(range, quantile, summary, fivenum)) {  
  print(f(y))  
}
```

[1] -8.304924 11.022464

##	0%	25%	50%	75%	100%	
##	-8.304924	-3.058389	1.325795	2.913598	11.022464	
##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	-8.3049	-3.0584	1.3258	0.4374	2.9136	11.0225
##	[1]	-8.304924	-3.309386	1.325795	2.951253	11.022464

Correlations

```
cor(x1, y) # get the correlation
```

```
## [1] 0.295735
```

Correlations: statistical test

```
cor.test(x1, y) # run a test

##
##  Pearson's product-moment correlation
##
## data:  x1 and y
## t = 1.3134, df = 18, p-value = 0.2055
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
##  -0.1688880  0.6528217
## sample estimates:
##          cor
## 0.295735
```


Correlations: using results

```
res = cor.test(x1, y) # save the result
names(res) # see result structure

## [1] "statistic"    "parameter"    "p.value"      "estimate"
## [6] "alternative"  "method"       "data.name"    "conf.int"

R = res['statistic'] # extract values...
p = res['p.value'] # ...from the result
```

Null hypothesis significance testing

Independent 2-sample 2-sided t-test

Test whether difference in means is different from 0

```
t.test(x1, y)
```

```
##
```

```
##  Welch Two Sample t-test
```

```
##
```

```
## data:  x1 and y
```

```
## t = 0.61544, df = 19.607, p-value = 0.5453
```

```
## alternative hypothesis: true difference in means is not equal to 0
```

```
## 95 percent confidence interval:
```

```
##  -1.520858  2.791563
```

```
## sample estimates:
```

```
## mean of x mean of y
```

```
## 1.0727840 0.4374315
```

More t-tests

The `t.test` function also has options for:

- 1-sided tests
- paired tests
- 1-sample tests.

Regression models

The `lm` (linear model) function and variants are used for regression.

```
df = tibble(x1, x2, y)
head(df)
```

```
## # A tibble: 6 x 3
##       x1       x2       y
##   <dbl> <dbl> <dbl>
## 1  1.94  -0.220  1.07
## 2  1.62   0.702  5.72
## 3  0.514  0.505  2.88
## 4  1.27   0.384  4.29
## 5  0.475 -0.964 -4.36
## 6  0.132 -0.512 -2.81
```

R provides a special formula syntax involving the tilde `~`. It's used to specify a regression model. The left-hand side is the dependent variable, `y`. The right-hand side gives the independent variables, interactions, and transformations. So, `~` means something like “is modelled as”.

`y ~ x1 + x2`

This says: run the formula $y = a + b_1x_1 + b_2x_2$

Using a formula in a regression

```
res <- lm(y ~ x1 + x2, data=df)
summary(res) # show results
```

```
##
## Call:
## lm(formula = y ~ x1 + x2, data = df)
##
```

```
## Residuals:
```

```
##      Min       1Q   Median       3Q      Max
## -3.0111 -0.6745  0.2500  0.6999  2.7058
```

```
##
```

```
## Coefficients:
```

```
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -0.7509     0.7119  -1.055   0.3063
## x1             1.2890     0.5867   2.197   0.0422 *
## x2             4.5939     0.3722  12.343 6.53e-10 ***
```

```
## ---
```

Formulas with interaction

If we changed + to *, we would add the interaction effect, ie we would run the formula

$$y = a + b_1x_1 + b_2x_2 + b_{12}x_1x_2$$

Use `?formula` for more on this special syntax.

Formulas with interaction

```
res <- lm(y ~ x1 * x2, data=df)
summary(res) # show results
```

```
##
## Call:
## lm(formula = y ~ x1 * x2, data = df)
```

```
##
## Residuals:
```

	Min	1Q	Median	3Q	Max
##	-2.6828	-0.4217	0.2195	0.7291	2.2952

```
##
## Coefficients:
```

	Estimate	Std. Error	t value	Pr(> t)	
## (Intercept)	-0.6443	0.7051	-0.914	0.3744	
## x1	1.2612	0.5774	2.184	0.0442	*
## x2	5.7138	0.9635	5.930	2.11e-05	***
## x1:x2	-1.0724	0.8535	-1.257	0.2270	

Formulas with transformation

We could also use transformations. For example:

```
res <- lm(y ~ x1 + log(x2), data=df)
```

```
## Warning in log(x2): NaNs produced
```

```
summary(res) # show results
```

```
##
```

```
## Call:
```

```
## lm(formula = y ~ x1 + log(x2), data = df)
```

```
##
```

```
## Residuals:
```

```
##      Min       1Q   Median       3Q      Max
```

```
## -2.9823 -1.7985 -0.0445  0.8793  4.2214
```

```
##
```

```
## Coefficients:
```

```
##              Estimate Std. Error t value Pr(>|t|)
```

```
## (Intercept)   4.4794      2.0797   2.154  0.0747 .
```

```
## x1           0.8525      1.6607   0.514  0.6257
```

One-way analysis of variance (ANOVA)

Like t-test for multiple groups, again using a formula.

```
res = aov(height ~ gender * species, data=dplyr::starwars)
summary(res)
```

```
##              Df Sum Sq Mean Sq F value    Pr(>F)
## gender          3    1674    558.1     8.421 0.000254 ***
## species        34   73457   2160.5    32.599 < 2e-16 ***
## gender:species   2     196     97.9     1.477 0.242539
## Residuals       34    2253     66.3
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 13 observations deleted due to missingness
```

Beyond Base R: the caret package

- k-nearest neighbours
- Linear regression
- Support vector machines
- Classification/regression trees
- Perceptrons
- Ensembles, including forests, bagging, boosting

<https://topepo.github.io/caret>

The caret package

The main Python competitor is `scikit-learn` which we will study later.

We won't go into detail on ML algorithms in this class.

Further reading

- <https://www.statmethods.net/stats/ttest.html>
- <https://www.statmethods.net/stats/regression.html>
- <https://www.statmethods.net/stats/anova.html>

Exercises

- 1 In the `mpg` dataset (part of the tidyverse), calculate the mean and standard deviation of the highway fuel efficiency.
- 2 Using `group_by`, calculate the mean and standard deviation of the highway fuel efficiency per manufacturer.
- 3 Calculate the correlation between highway fuel efficiency and engine size.
- 4 What was the average highway fuel efficiency in 1999 and in 2008?
- 5 Carry out a two-sample independent t-test between highway fuel efficiency in 1999 and 2008 and interpret the result.
- 6 Carry out a regression on highway fuel efficiency by displacement.

Solution 1

```
library(tidyverse)
```

```
mean(mpg$hwy)
```

```
## [1] 23.44017
```

```
sd(mpg$hwy)
```

```
## [1] 5.954643
```


Solution 2

```
mpg %>% group_by(manufacturer) %>%  
  summarise(mean=mean(hwy), sd=sd(hwy))
```

```
## # A tibble: 15 x 3  
##   manufacturer mean    sd  
##   <chr>      <dbl> <dbl>  
## 1 audi       26.4  2.18  
## 2 chevrolet  21.9  5.11  
## 3 dodge      17.9  3.57  
## 4 ford       19.4  3.33  
## 5 honda      32.6  2.55  
## 6 hyundai    26.9  2.18  
## 7 jeep       17.6  3.25  
## 8 land rover 16.5  1.73  
## 9 lincoln    17    1  
## 10 mercury   18    1.15  
## 11 nissan     24.6  5.09
```

Solution 3

```
cor(mpg$hwy, mpg$displ)
```

```
## [1] -0.76602
```

Solution 4

```
mpg %>% group_by(year) %>%  
  summarise(mean=mean(hwy), sd=sd(hwy))
```

```
## # A tibble: 2 x 3  
##   year mean    sd  
##   <int> <dbl> <dbl>  
## 1  1999  23.4  6.08  
## 2  2008  23.5  5.85
```

Solution 5

```
mpg1999 <- mpg %>% filter(year == 1999)
mpg2008 <- mpg %>% filter(year == 2008)
t.test(mpg1999$hwy, mpg2008$hwy)
```

```
##
##  Welch Two Sample t-test
##
## data:  mpg1999$hwy and mpg2008$hwy
## t = -0.032864, df = 231.64, p-value = 0.9738
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
##  -1.562854  1.511572
## sample estimates:
## mean of x mean of y
##  23.42735  23.45299
```

Solution 6

```
res = lm(hwy ~ displ, data=mpg)
summary(res)

##
## Call:
## lm(formula = hwy ~ displ, data = mpg)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -7.1039 -2.1646 -0.2242  2.0589 15.0105
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  35.6977     0.7204   49.55  <2e-16 ***
## displ       -3.5306     0.1945  -18.15  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```