

Tables in R – A quick practical overview

<preliminary blueprint version>

by Andri Signorell

Helsana Versicherungen AG, Health Sciences, Zurich

HWZ University of Applied Sciences in Business Administration, Zurich

andri@signorell.net

October, 24th, 2015

Tabulating data is both, trivial and complicated. After all it is just about counting data. But the underlying data structures are diverse and technically abstract, especially when there are more than two dimensions involved. Thus there are many functions to handle and process tables in the respective representation, which makes the situation somewhat confusing. There are some gaps in base R function list that are filled by DescTools.

This document aims to briefly summarise, how to create, manipulate and describe count data in tables. Some examples from the SAS-documentation FREQ are reproduced.

1	Starting Point.....	2
2	Create Table	3
2.1	Creating from the scratch.....	3
2.2	Building categories from a numeric variable.....	4
2.3	Expanding.....	5
2.4	SAS datalines	5
3	Tabulate	5
4	Reorganize.....	7
5	Aggregate	9
6	Append.....	10
7	Convert.....	10
8	Print and Format.....	11
9	Export.....	12
10	Plot.....	13
11	Save.....	14
12	Descriptions, Statistics and Tests	14
13	Cases	19
13.1	Eye colour - Binomial Proportions for One-Way Frequency Tables	19
13.2	Cochran-Armitage Trend Test.....	20
13.3	Heart - 2x2-Table.....	20
13.4	Ophtalmological Test – Confusion Matrix	23
13.5	Skin - Agreement Study	24
13.6	Migraine - Statistics for a Stratified 2x2-Table	25
14	References.....	29

Note: For all the examples in this document, `library(DescTools)` must be declared.

1 Starting Point

The analysis of categorical data usually starts with tables. In R we have a comprehensive but not complete toolset to work with tables of two and more dimensions. Thus there's room for extensions with functions useful in the analyst's daily life. DescTools contains quite a few of such, which are described in this document.

The first question is how categorical data is technically organised. Normally it will be given in one of the following three data structures.

A) Single case

The raw data in form of a `data.frame` (or a matrix), each row contains one case, here one person:

```
head(d.col, 20)
```

	Hair	Eye	Sex
1	Black	Brown	Male
2	Black	Brown	Male
3	Black	Brown	Male
4	Black	Brown	Male
5	Black	Brown	Male
6	Black	Brown	Male
7	Black	Brown	Male
...			
32	Black	Brown	Male
33	Brown	Brown	Male

...

B) Frequency

Unique combinations of factors extended with their counts, often called weights (column "Freq"):

```
d.weight
```

	Hair	Eye	Sex	Freq
1	Black	Brown	Male	32
2	Brown	Brown	Male	53
3	Red	Brown	Male	10
4	Blond	Brown	Male	3
5	Black	Blue	Male	11
6	Brown	Blue	Male	50
7	Red	Blue	Male	10
8	Blond	Blue	Male	30
9	Black	Hazel	Male	10
10	Brown	Hazel	Male	25

...

C) Table

A multidimensional table (or an array, matrix):

```
tab
```

```
, , Sex = Male
      Eye
Hair   Brown Blue Hazel Green
Black   32   11   10     3
Brown   53   50   25   15
Red     10   10    7    7
Blond    3   30    5    8

, , Sex = Female
      Eye
Hair   Brown Blue Hazel Green
...
```

...

Either we have the raw data organised case-by-case in a `data.frame` (case A). Then a contingency table can be built by tabulating the data. There are several commands for this described in chapter "Tabulate".

Or the data are given as a combination of factor levels and one count variable (typically organised as a `data.frame` too) (case B). The first line in this representation means, that we have 32 men with black hair and brown eyes in our sample. This corresponds to the cell [1, 1, 1] in the representation C). In representation A we have 32 rows with the same content Black/Brown/Male.

How to create directly such a structure is described in "Expanding". There are functions to convert this structure to a table or to recreate the raw dataset. This is detailed in the chapter "Convert".

When the data are given directly as a table (case C), there are again several ways how to enter that into R. This is the content of the first chapter "Create tables".

How to process tables is described in the chapters "Reorganize", "Aggregate", "Append", "Convert".

Usually B) will be the most economic representation of frequency data whereas the case-by-case form in A) is the least (provided the data set is purely categorical). The built-in data sets from the R base system that are purely categorical usually come in the form of tables (C).

2 Create Table

2.1 Creating from the scratch

There are several ways to enter contingency table data into R. Let's illustrate here some approaches with a table concerning party affiliation by gender:

Table 2.1 Tabulating Party versus Gender, Agresti (2007) p. 39

Gender	Party		
	Democrat	Independent	Republican
M	762	327	468
F	484	239	477

The first approach uses the function `rbind` and builds a matrix row by row. The `as.table()` function lets R know that the matrix represents a contingency table of counts.:

```
tab <- as.table(rbind(c(762, 327, 468), c(484, 239, 477)))
dimnames(tab) <- list(gender = c("M", "F"),
                      party = c("Democrat", "Independent", "Republican"))
tab

##      party
## gender Democrat Independent Republican
##      M      762      327      468
##      F      484      239      477
```

rbind
as.table

The exactly same result can be created by the second approach, using the function `matrix`. Note that, by default, `matrix()` uses the elements supplied by columns in the result, unless you specify `byrow=TRUE`.

```
as.table(matrix(c(762, 327, 468, 484, 239, 477), nrow=2, byrow=TRUE,
                dimnames=list(gender= c("M", "F"),
                              party = c("Democrat", "Independent", "Republican"))))
```

matrix

The third way uses `TextToTable` to convert a text to a table. Within this function `read.table` is used to enter the data and to convert the data.frame to a table. `header=TRUE` will take the names of the variables from its first line.

TextToTable

The column names and row names will automatically be chosen, if the first row contains one fewer field than the number of columns. The dimension names can be provided with the specific argument:

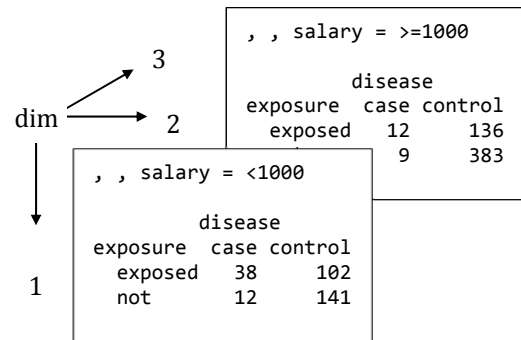
```
txt <- "
  Democrat, Independent, Republican
M, 762, 327, 468
F, 484, 239, 477"

TextToTable(txt, sep=",", dimnames=c("gender", "party"))
```

Higher dimensional arrays can be defined with the function `array` by using the argument `dim`:

array

```
salary <- array(
  c(38, 12, 102, 141, 12, 9, 136, 383),
  dim=c(2, 2, 2),
  dimnames=list(exposure = c("exposed", "not"),
                disease = c("case", "control"),
                salary = c("<1000", ">=1000"))
)
```



Note how the dimensions are organised:
The first dimension corresponds to the rows,
the second to the columns, the third to the depth, and so on.

Figure 2.1 3-dimensional table

Higher dimensional tables condensed in flat tables with more than one column, resp. row variable, can be created from the appropriate text chunk by means of the base function `read.ftable`. (Beware not to insert spaces at the beginning of the lines.)

read.ftable

```
txt <-
"
      Sex  Male                Female
      Eye Brown Blue Hazel Green  Brown Blue Hazel Green
Hair
Black      32   11    10     3    36   9    5    2
Brown      53   50    25    15    66   34   29   14
Red        10   10     7     7    16   7    7    7
Blond       3   30     5     8     4   64   5    8
"
tab <- as.table(read.ftable(textConnection(txt)))
```

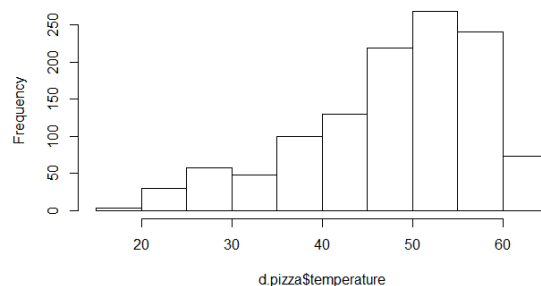
2.2 Building categories from a numeric variable

When a numeric variable has to be cut into intervals there's the function `hist()` for creating a histogram. The DescTools function `Freq()` is designed to give the numeric representation of a histogram. It displays the frequencies and the percentages of a binned variable with the same default logic as `hist()`. The single and cumulative frequencies values are reported.

```
Freq(d.pizza$temperature)

##      level freq  perc cumfreq cumperc
## 1  [15,20]    3 0.003     3  0.003
## 2  (20,25]   30 0.026    33  0.028
## 3  (25,30]   58 0.050    91  0.078
## 4  (30,35]   48 0.041   139  0.119
## 5  (35,40]  100 0.085   239  0.204
## 6  (40,45]  130 0.111   369  0.315
## 7  (45,50]  219 0.187   588  0.503
## 8  (50,55]  268 0.229   856  0.732
## 9  (55,60]  241 0.206  1097  0.938
## 10 (60,65]   73 0.062  1170  1.000

hist(d.pizza$temperature)
```



Freq

hist

Figure 2.2 Histogram of a numeric variable.

2.3 Expanding

For small frequency tables, it is often convenient to enter them in frequency form using `expand.grid()` for the factors and `c()` to list the counts in a vector.

```
tab <- data.frame(expand.grid(
  Hair = c("Black", "Brown", "Red", "Blond"),
  Eye  = c("Brown", "Blue", "Hazel", "Green"),
  Sex  = c("Male", "Female")),
  count = c(32,53,10,3,11,50,10,30,10,25,7,5,3,15,7,8,
            36,66,16,4,9,34,7,64,5,29,7,5,2,14,7,8) )
```

expand.grid

`expand.grid` will create all the interactions between the given factors. `data.frame` will bind them with the count variable, denominating the number of observations. This will be a type B representation of count data, which can be converted with `xtabs` to a table. (See Chapter “Convert”)

2.4 SAS datalines

Longstanding predominance of SAS entails, that small data tables in examples and documents are often reported in the SAS datalines format. Creating a table based on this in R is not straight forward, as there might be more than one case per row (as in the example below).

The function `ParseSASDatalines` parses the syntax and creates a table named after the *data* statement, using given column names (specified by the keyword *input*).

```
ParseSASDatalines("
  data SummerSchool;
  input Gender $ Internship $ Enrollment $ Count @@;
  datalines;
  boys  yes yes 35  boys  yes no 29
  boys  no  yes 14  boys  no  no 27
  girls yes yes 32  girls yes no 10
  girls no  yes 53  girls no  no 23
;")
```

**ParseSAS-
Datalines**

The command above will directly (and silently) create a new data object named *SummerSchool* in the GlobalEnvironment.

3 Tabulate

The built-in data set *HairEyeColor* has the class *table*. Let's turn this table into a case-by-case data frame as a base for the subsequent analysis. `Untable` does this job.

```
d.col <- Untable(HairEyeColor)
head(d.col, 3)
```

Untable

```
## Hair  Eye  Sex
## 1 Black Brown Male
## 2 Black Brown Male
## 3 Black Brown Male
```

From here we can start tabulating again. The simplest case is to tabulate a single vector. The function `table` yields the absolute frequencies and `prop.table` the proportions:

```
table(d.col$Hair)                                prop.table(table(d.col$Hair))

## Black Brown  Red Blond                        ##   Black   Brown   Red   Blond
##   108   286   71  127                        ## 0.1824324 0.4831081 0.1199324 0.2145270
```

**table
prop.table**

A combination of both extended with the cumulative sums for both, absolute and relative frequencies, can be produced by `Freq` (here ordered by decreasing frequency):

```
Freq(d.col$Hair, ord="desc")
```

Freq

```
##   level freq  perc cumfreq cumperc
## 1 Brown  286 0.483    286   0.483
## 2 Blond  127 0.215    413   0.698
## 3 Black  108 0.182    521   0.880
## 4 Red    71 0.120    592   1.000
```

By means of the `table` function we can produce multidimensional contingency tables (aka. crosstabs) as well. We use the command `with` here, so we can avoid having to qualify every column name with the name of the data.frame (which makes the code more readable).

```
with(d.col, table(Hair, Eye))
```

```
##           Eye
## Hair      Brown Blue Hazel Green
## Black      68   20   15     5
## Brown     119   84   54    29
## Red        26   17   14    14
## Blond       7   94   10    16
```

The first entered variable will be the row variable, the second one the column variable.

Missing values are ignored by default. In order to include NA as a category in counts, use the option `useNA="always"`.

A relative frequency table can be produced using the function `prop.table`, which takes a table object as argument:

```
with(d.col, prop.table(table(Hair, Eye), margins=NULL))
```

```
##           Eye
## Hair      Brown      Blue      Hazel      Green
## Black 0.114864865 0.033783784 0.025337838 0.008445946
## Brown 0.201013514 0.141891892 0.091216216 0.048986486
## Red   0.043918919 0.028716216 0.023648649 0.023648649
## Blond 0.011824324 0.158783784 0.016891892 0.027027027
```

The function `PercTable` combines that and allows adding marginal sums in one step:

```
PercTable(Hair ~ Eye, data=d.col, rfrq="111", margins=c(1,2))
```

PercTable

```
##           Eye  Brown   Blue  Hazel  Green   Sum
## Hair
## Black freq      68    20    15     5    108
##      perc     .115   .034   .025   .008   .182
##      p.row     .630   .185   .139   .046     .
##      p.col     .309   .093   .161   .078     .
## Brown freq     119    84    54    29    286
##      perc     .201   .142   .091   .049   .483
##      p.row     .416   .294   .189   .101     .
##      p.col     .541   .391   .581   .453     .
## Red  freq      26    17    14    14     71
##      perc     .044   .029   .024   .024   .120
##      p.row     .366   .239   .197   .197     .
##      p.col     .118   .079   .151   .219     .
## Blond freq       7    94    10    16    127
##      perc     .012   .159   .017   .027   .215
##      p.row     .055   .740   .079   .126     .
##      p.col     .032   .437   .108   .250     .
##
## Sum  freq      220    215    93    64    592
##      perc     .372   .363   .157   .108   1.000
##      p.row     .   .   .   .   .
##      p.col     .   .   .   .   .
```

There are more options, as expected values or standard residuals, which can optionally be integrated.

The marginal tables can be produced by R base function `margin.table` or by the somewhat extended function `Margins` in `DescTools`:

```
Margins(tab, ord="desc")
```

Margins

```
## $Hair
##   level freq  perc cumfreq cumperc
## 1 Brown  286 0.483    286    0.483
## 2 Blond  127 0.215    413    0.698
## 3 Black  108 0.182    521    0.880
## 4  Red   71 0.120    592    1.000

## $Eye
##   level freq  perc cumfreq cumperc
## 1 Brown  220 0.372    220    0.372
## 2  Blue  215 0.363    435    0.735
## 3 Hazel   93 0.157    528    0.892
## 4 Green   64 0.108    592    1.000
```

`table` does not come with a formula interface, but the `xtabs` function does. This allows us to create multidimensional crosstabulations using formula style input. The result is a contingency table in array format, whose dimensions are determined by the terms on the right side of the formula.

4 Reorganize

Say we created a three dimensional table with *Hair*, *Eye* and *Sex* as variables and typically got a 3-dim array as result. This will be displayed as:

```
(tab <- with(d.col, table(Hair, Eye, Sex)))
```

```
## , , Sex = Male
##
##      Eye
## Hair  Brown Blue Hazel Green
## Black    32   11   10    3
## Brown    53   50   25   15
## Red      10   10    7    7
## Blond     3   30    5    8

## , , Sex = Female
##
##      Eye
## Hair  Brown Blue Hazel Green
## Black    36    9    5    2
## Brown    66   34   29   14
## Red      16    7    7    7
## Blond     4   64    5    8
##
```

To combine this multidimensional structure into a flat table while preserving all the details, there's the function `ftable`. The variables to be placed in the rows can be defined by the argument `row.vars`, which can be a vector (denoting multiple dimensions) containing the dimension numbers or the names, if there are any defined.

So to put *Eye* (variable 2) and *Sex* (variable 3) in the rows and *Hair* as column variable, we can use both, subscripts or dimension names, writing

```
ftable(tab, row.vars = c(2, 3))
ftable(tab, row.vars = c("Eye", "Sex"))

##           Hair Black Brown Red Blond
## Eye  Sex
## Brown Male          32   53  10    3
##        Female        36   66  16    4
## Blue  Male          11   50  10   30
##        Female         9   34   7   64
## Hazel Male          10   25   7    5
##        Female         5   29   7    5
## Green Male           3   15   7    8
##        Female         2   14   7    8
```

ftable

The tab, as we constructed it, has the *Hair* as rows (1), the *Eye* as columns (2), and the *Sex* as third dimension (3) defined. The dimensions and dimension names follow the defined order:

```
dimnames(tab)

## $Hair
## [1] "Black" "Brown" "Red"   "Blond"
##
## $Eye
## [1] "Brown" "Blue"  "Hazel" "Green"
##
## $Sex
## [1] "Male"  "Female"
```

If we have to change the order of the dimensions, we can make use of the base-R function `aperm`. Let's say we wanted *Eye* as row variable and *Sex* a column variable and consequently *Hair* as 3th variable, we can tell `aperm` to set dimension 2 on the first position, 3 on the second and 1 on the third position. So we get:

aperm

```
aperm(tab, c(2,3,1))

##, , Hair = Black
##
##      Sex
## Eye   Male Female
## Brown  32      36
## Blue   11       9
## Hazel  10       5
## Green   3       2
##

##, , Hair = Brown
##
##      Sex
## Eye   Male Female
## Brown  53      66
## Blue   50      34
##
##
...

```

The following would by the way not work:

```
tab["Eye", "Sex", "Hair"]
Error in tab["Eye", "Sex", "Hair"] : subscript out of bounds
```

To reorder the sequence of the levels (within a dimension) in our table, we could use `reorder.factor`. Say we would like to have the sequence *Blue, Green, Hazel, Brown* for the *Eye* colour. Of course, when having the raw data, we would use

```
factor(d.col$Eye, levels=c("Blue", "Green", "Hazel", "Brown"))
```

and any table afterwards would inherit this level order. But how can we change this in an already created table? The answer is obvious (but may yet be unexpected in this context): Use the subscript! This works with the level names as well as with the index positions.

```
tab[, c("Blue", "Green", "Hazel", "Brown"), ]

##, , Sex = Male
##
##      Eye
## Hair   Blue Green Hazel Brown
## Black  11     3   10    32
## Brown  50    15   25    53
## Red    10     7    7     10
## Blond  30     8    5     3
##

##, , Sex = Female
##
##      Eye
## Hair   Blue Green Hazel Brown
## Black   9     2    5    36
## Brown  34    14   29    66
## Red     7     7    7    16
## Blond  64     8    5     4

```


For simply reversing the levels there's the function `Rev`, which has a table interface implemented. The function accepts a `margins` argument, defining the dimensions whose levels should be reversed. Compare the reversed levels of *Hair* and *Sex*:

Rev

```
tab                                Rev(tab, margin = c(1, 3))

## , , Sex = Male                  ## , , Sex = Female
##                                ##
##      Eye                        ##      Eye
## Hair   Brown Blue Hazel Green  ## Hair   Brown Blue Hazel Green
## Black   32  11  10   3         ## Blond   4  64   5   8
## Brown   53  50  25  15         ## Red     16  7   7   7
## Red     10  10   7   7         ## Brown   66  34  29  14
## Blond   3  30   5   8         ## Black   36  9   5   2
##                                ##
## , , Sex = Female                ## , , Sex = Male
##                                ##
##      Eye                        ##      Eye
## Hair   Brown Blue Hazel Green  ## Hair   Brown Blue Hazel Green
## Black   36  9   5   2         ## Blond   3  30   5   8
## Brown   66  34  29  14         ## Red     10  10   7   7
## Red     16  7   7   7         ## Brown   53  50  25  15
## Blond   4  64   5   8         ## Black   32  11  10   3
```

Renaming level names can be achieved by refining the dimension names.

```
dimnames(tab)$Sex <- c("men", "women")
```

5 Aggregate

Sometimes we might want to aggregate an existing table along one or several dimensions.. Say we'd like to get rid of the *Hair* dimension, but retain all the frequency information for the other dimensions. For this we can use `apply` as we would in the case of a matrix. The function takes as well vectors for the margins. The order of the subscripts specified in the `apply` statement determines the order of the subscripts in the result.

So if we sum up all cases along the 1st dimension (*Hair*) and retain the other two (2, 3) we would get:

apply

```
apply(tab, c(2,3), sum)           apply(tab, 1, sum)

##      Sex                        ## Black Brown  Red Blond
## Eye   Male Female              ##  108  286   71  127
## Brown  98   122
## Blue  101   114
## Hazel  47    46
## Green  33    31
```

This works with the dimension names too: `apply(tab, c("Eye", "Sex"))` will deliver the same result.

Single margins could be calculated correspondingly, as demonstrated above.

If `tab` was created with `xtabs`, it can be aggregated directly by using the formula interface, which typically is clearer and more readable.

xtabs

```
xtab <- xtabs(~., d.col)
xtabs(Freq ~ Eye + Sex, xtab)

##      Sex
## Eye   Male Female
## Brown  98   122
## Blue  101   114
## Hazel  47    46
## Green  33    31
```

If we want to combine some levels, we can with `CollapseTable`. Say we want to fuse brown and hazel eyes to a new category *Browny*, as well as just having two groups of *Hair*, namely *Dark* and *Fair*:

```
CollapseTable(tab, Eye=c("Brown", "Blue", "Brown", "Green"),
              Hair=c("Dark", "Dark", "Fair", "Fair"))
```

CollapseTable

```
## , , Sex = Male
##
##      Eye
## Hair   Brown Blue Green
## Dark   120   61   18
## Fair    25   40   15

## , , Sex = Female
##
##      Eye
## Hair   Brown Blue Green
## Dark   136   43   16
## Fair    32   71   15
```

6 Append

Sometimes we need to paste tables together, for instance when two tables of the same dimension should be put together to a 3-dimensional array. In contrast to the 2-dimensional case, where the functions `rbind` and `cbind` exist, base R does not contain a respective function for higher dimensional tables. In `DescTools` there's the function `Abind` included for this purpose (indeed borrowed from the `abind` package).

```
a <- HairEyeColor[,1] # male table
b <- HairEyeColor[,2] # female table
```

```
Abind(Male=a, Female=b, along=3)
```

Abind

```
## , , Male
##
##      Brown Blue Hazel Green
## Black   32   11   10    3
## Brown   53   50   25   15
## Red     10   10    7    7
## Blond    3   30    5    8
##

## , , Female
##
##      Brown Blue Hazel Green
## Black   36    9    5    2
## Brown   66   34   29   14
## Red     16    7    7    7
## Blond    4   64    5    8
```

The first step separates the table for males from the females. `Abind` reverses this step and binds the two tables together again. This can happen along all possible dimensions.

In the example above a new dimension is introduced by setting `along = 3`.

`Abind(a,b,along=2)` would bind the tables by columns (as `cbind` does), whereas `Abind(a,b,along=1)` would give the same result as `rbind(a,b)`.

7 Convert

Time and again newbies wonder how to convert tables from one to the other form. Base R comprises most of the required functions, but not quite all.

Let's say we have the three forms of table given as:

```
d.col <- Untable(HairEyeColor) # case-by-case, A)
d.weight <- as.data.frame(HairEyeColor) # frequency, B)
tab <- HairEyeColor # table, C)
```

The conversions can be made as follows.

A) → B) 1) `as.data.frame(table(d.col))`
 2) `aggregate(rep(1, nrow(d.col)),`
 `by=d.col, FUN=length)`

A) ← B) `Untable(d.weight)`

A) → C) `table(d.col)`

A) ← C) `Untable(tab)`

 B) → C) `xtabs(Freq ~ ., d.weight)`

 B) ← C) `as.data.frame(tab)`

This is actually A) to C) to B)!

Solution 2) will yield the nonzero entries only.

`library(DescTools)`

If `tab` is defined as matrix, `as.data.frame` has to be specified explicitly as `as.data.frame.table!`

as.data.frame

aggregate

Untable

xtabs

The conversion of an `xtabs` object to a matrix would normally not be carried out in base R. The class would remain (“`xtabs`”, “`table`”) after calling the `as.matrix` function. All attributes won’t be touched as well.

```
str(as.matrix(htab))
##  xtabs [1:2, 1:7] 119 1070 16 60 12 14 7 4 3 0 ...
##  - attr(*, "dimnames")=List of 2
##    ..$ race : chr [1:2] "Black" "White"
##    ..$ nvics: chr [1:7] "0" "1" "2" "3" ...
##  - attr(*, "class")= chr [1:2] "xtabs" "table"
##  - attr(*, "call")= language xtabs(formula = freq ~ race + nvics,
##                                data = homicide)
```

`DescTools` will add a `xtabs` interface for `as.matrix` such, that the class and call attributes will be adapted.

```
library(DescTools)
## Loading required package: manipulate

str(as.matrix(htab))
##  num [1:2, 1:7] 119 1070 16 60 12 14 7 4 3 0 ...
##  - attr(*, "dimnames")=List of 2
##    ..$ race : chr [1:2] "Black" "White"
##    ..$ nvics: chr [1:7] "0" "1" "2" "3" ...
```

8 Print and Format

All table connected classes have their print methods which do not call for any further explanation. There are several approaches out there, how to turn tables into XML, HTML or LATEX. `DescTools` contains two functions for sending tables to MS-Word. `WrdTable` would create the table in Word and transfer the cell information appropriately.

Let’s create an artificial table, with one cell being 0 and one being NA. Then we format the counts with a big.mark and set 0 digits. The zero values should be expressed as “-” and the NAs as “missing”. Finally all should be aligned to the right.

```
(tab <- as.table(matrix(c(2000, 0, 34, NA), nrow=2)))
##      A      B
## A 2000    34
## B      0
```

```
tab[] <- Format(tab, big.mark = "'", digits=0, zero.form="-", na.form="Missing")
tab[] <- StrAlign(tab, "\\$")
tab
```

Format

```
##           A           B
## A    2'000           34
## B           - Missing
```

The counts and percentages in `PercTable` can be formatted by setting the options `fmt.abs` and `fmt.per`. The percentages are formatted as `.000` and the counts with a space for `big.mark`.

```
options(fmt.abs=structure(list(digits=1, big.mark=" "), class="fmt"))
options(fmt.per=structure(list(digits=3, leading="drop"), class="fmt"))
PercTable(tab)
```

```
##           A           B
##
## A freq  2 000.0      34.0
##   perc    .720     .012
##
## B freq    0.0     745.0
##   perc    .000     .268
```

Note that by applying formats to the cells, the numeric values turn to strings and cannot be subsequently used for further calculating.

`FixToTab` is trying to chop the fixed font output of a table given as text to a tab delimited table.

9 Export

`DescTools` contains functions for exporting tables to Word or Excel. Exporting to Excel would at least export the values. (This can – and should – be improved..)

```
tab <- ftable(HairEyeColor, col.vars = c("Sex", "Hair"))
XLView(tab)
```

XLView

	A	B	C	D	E	F	G	H	I	J
1		V1	V2	V3	V4	V5	V6	V7	V8	
2	1	32	53	10	3	36	66	16	4	
3	2	11	50	10	30	9	34	7	64	
4	3	10	25	7	5	5	29	7	5	
5	4	3	15	7	8	2	14	7	8	
6										

Figure 9.1 Excel sheet containing exported table from R.

The Word-Interface is already somewhat more elaborated (but still unsatisfactory):

```
WrdTable(tab, row.names = FALSE, wrd=GetNewWrd())
```

WrdTab

	Sex	Male				Female			
	Hair	Black	Brown	Red	Blond	Black	Brown	Red	Blond
Eye									
Brown		32	53	10	3	36	66	16	4
Blue		11	50	10	30	9	34	7	64
Hazel		10	25	7	5	5	29	7	5
Green		3	15	7	8	2	14	7	8

10 Plot

The usual representation of a table is a mosaicplot. Such a plot will display the conditional frequencies in two directions. (Note that the important encoding is length.)

```
tab <- as.table(apply(HairEyeColor, c(1,2), sum))
tab <- tab[,c("Brown", "Hazel", "Green", "Blue")]
cols <- SetAlpha(c("sienna4", "burlywood", "chartreuse3", "slategray1"), 0.6)

PlotMosaic(tab, col=cols, main = "Hair ~ Eye")
```

This will display the following fact:

PlotMosaic

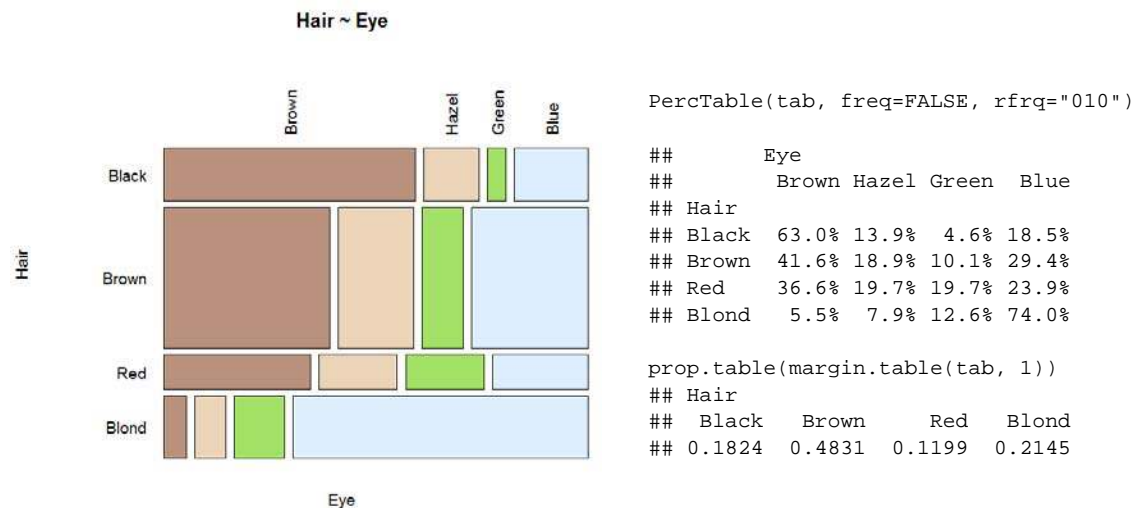


Figure 10.1 Mosaicplot of Hair colour ~ Eye colour.

The plot makes the story quite visible! About half of the sample has brown hair, red is the less frequent hair colour observed (~5-10%). Within the black haired people more than 50% have brown eyes. Blond people tend to have blue eyes. The percentage of green eyed people is biggest within red haired guys, but with 20% not as pronounced as maybe expected. And so on.

The mosaicplot has an order. At first the hair colour is split and afterwards, within the single hair colour, the eye colour. This corresponds to a relationship Hair ~ Eye. If the inverse relation is the interesting one, the table can simply be transposed. This side of the coin then looks like:

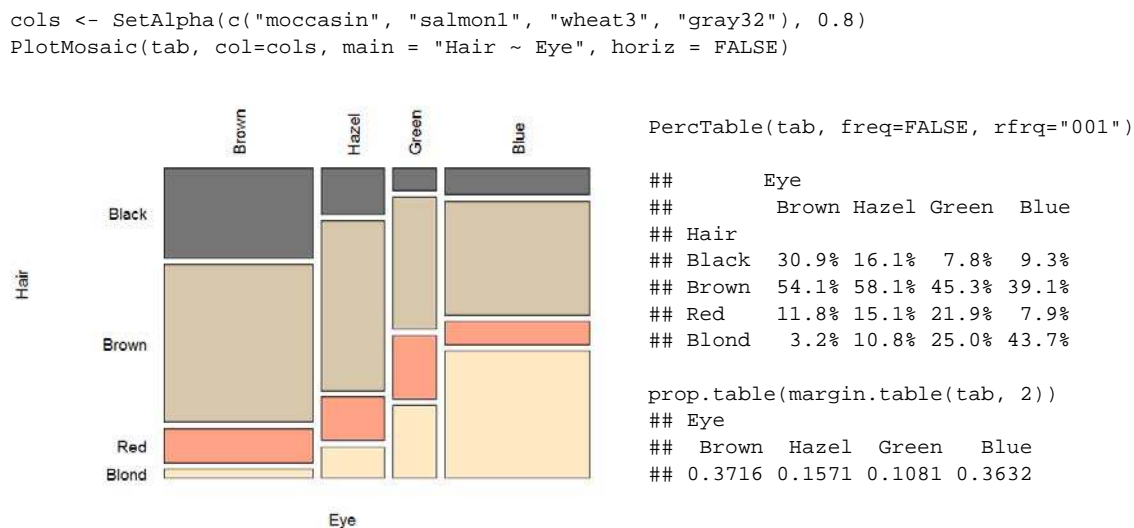
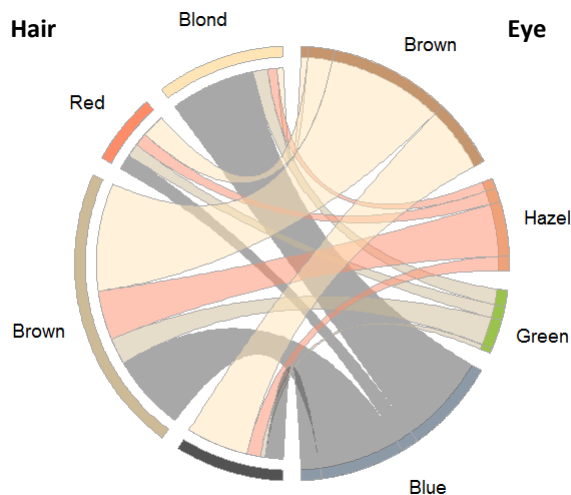


Figure 10.2 Mosaicplot of Eye colour ~ Hair colour.

Another – rather new – idea is to describe proportions in circles. It emphasises the association structure of the data. The left side of the circle represents the rows, say the hair colour, and the right one the columns, thus the eye colour. The advantage is that we see both marginal densities in the plot.



```
cols <- c("moccasin", "salmon1",
          "wheat3", "gray32",
          "slategray1", "chartreuse3",
          "burlywood", "sienna4")

PlotCirc(t(tab), acol=cols)
```

PlotCirc

Figure 10.3 Circular plot of HairEyeColor.

Looking at the blue eyes first of all we notice, that roughly a third of the sample has blue eyes. Within those, about 40% have blond hair, 10% red hair, 40% brown hair and again 10% black hair. When we follow the band from the blue eyed to the blond haired, we notice that blue eyed people form ~75% of the blond haired group. Obviously we see more (conditional) proportions in a circular plot than on a mosaic plot. A disadvantage is that angles are nowhere near as good to compare as the lengths in the mosaic.

11 Save

For saving the table, there's the usual R-base command:

```
save(tab, file = "HairEyeColor.rda")
```

save

12 Descriptions, Statistics and Tests

Let's create a 2-dimensional table and describe it with some bells and whistles. The argument `verbose = high` will maximize the volume of output:

```
# aggregate 3-d table to Eye and Hair colour only:
tab <- as.table(apply(HairEyeColor, c(2,3), sum))

# order the levels along colours:
tab <- tab[c("Brown", "Hazel", "Green", "Blue"),]
```

```
# describe the table
Desc(tab, verbose="high")
```

Desc

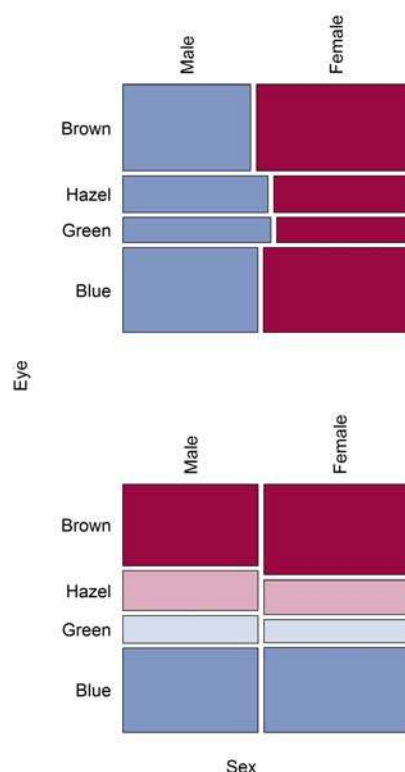
```
## Summary:
## n: 592, rows: 4, columns: 2
##
## Pearson's Chi-squared test:
##  X-squared = 1.5298, df = 3, p-value = 0.6754
## Pearson's Chi-squared test (cont. adj):
##  X-squared = 1.5298, df = 3, p-value = 0.6754
## Likelihood Ratio:
##  X-squared = 1.5294, df = 3, p-value = 0.6755
## Mantel-Haenszel Chi-squared:
##  X-squared = 0.2438, df = 1, p-value = 0.6214
```

The first line reports the total n in the table and the dimension, so we have 592 Persons in a table with 4 rows and two columns. Then several Chi-Square-tests are calculated. The null hypothesis is that the eye colour is not associated with the sex. The small value of the χ^2 -statistic, 1.5298, and the p-value of 0.6754 indicate that the null hypothesis can't be rejected at the 0.05 level of significance. Thus we would conclude that the observation does not indicate an association between eye colour and sex of the person.

The Pearson χ^2 -statistic involves the differences between the observed cell frequencies and the expected deviation-frequencies. Following a rule of thumb the expected frequency in every cell of the table should not be less than 5. R will print a message, if this condition is violated.

The continuity-adjusted χ^2 -test statistic consists of the Pearson χ^2 modified with an adjustment for continuity. As the sample size increases, the difference between the continuity-adjusted and Pearson χ^2 decreases. Thus in very large samples (as we have here) the two statistics are almost the same. This test statistic is also an alternative to Pearson's χ^2 if any of the expected values in a 2x2 table are less than 5. Some prefer to use the continuity-adjusted χ^2 -statistic when the sample size is small regardless of the expected values.

```
##           Sex
##           Male Female   Sum
## Eye
## Brown freq      98    122   220
##           perc    16.6%  20.6%  37.2%
##           p.row    44.5%  55.5%    .
##           p.col    35.1%  39.0%    .
##
## Hazel freq      47     46    93
##           perc     7.9%   7.8%  15.7%
##           p.row    50.5%  49.5%    .
##           p.col    16.8%  14.7%    .
##
## Green freq      33     31    64
##           perc     5.6%   5.2%  10.8%
##           p.row    51.6%  48.4%    .
##           p.col    11.8%   9.9%    .
##
## Blue  freq     101    114   215
##           perc    17.1%  19.3%  36.3%
##           p.row    47.0%  53.0%    .
##           p.col    36.2%  36.4%    .
##
## Sum   freq      279    313   592
##           perc    47.1%  52.9% 100.0%
##           p.row    .      .    .
##           p.col    .      .    .
```



The expected frequencies can be obtained by using the `expected` option on the `Desc` command (`Desc(tab, verbose="high", expected=TRUE)`). Additionally, the difference between the observed cell count and the expected cell count will be reported when using the `residuals=TRUE` and `stdres=TRUE` option for the standardized residuals (amount that each cell contributes to the value of the test statistic).

```
options(fmt.num=structure(list(digits=3), class="fmt"))
PercTable(tab, freq=TRUE, rfrq="000",
           expected=TRUE, residuals=TRUE)
```

		Sex	
		Male	Female
## Eye			
## Brown	freq	98	122
##	exp	103.682	116.318
##	res	-0.558	0.527
## Blue	freq	101	114
##	exp	101.326	113.674
##	res	-0.032	0.031

## Hazel	freq	47	46
##	exp	43.829	49.171
##	res	0.479	-0.452
## Green	freq	33	31
##	exp	30.162	33.838
##	res	0.517	-0.488

PercTable,
ExpFreq

This output shows the observed frequencies (`freq`), the expected values (`exp`) and the Pearson residuals (`res`), whose squared values are each cell's contribution to the χ^2 statistic. None of the expected values are less than 5, so we feel comfortable with the result of the Chi-Square test above.

The Likelihood Ratio χ^2 is asymptotically equivalent to the Pearson χ^2 (and Mantel-Haenszel χ^2) but not usually used when analyzing 2x2 tables. It is used in logistic regression and loglinear modeling which involves contingency tables.

The Mantel-Haenszel χ^2 is related to the Pearson χ^2 and, in the 2x2 case, as the sample size gets large these statistics converge. In the case of 2xC or Rx2 tables, if the variable with more than 2 categories is ordinal, the Mantel-Haenszel χ^2 is a test for trend while the Pearson χ^2 remains a general test for association.

When the `verbose` argument of the function `Desc` is set to "high", several statistics that describe the nominal and ordinal association between the two variables of the contingency table will be computed.

	estimate	lwr.ci	upr.ci
## Phi Coeff.	0.0508	-	-
## Contingency Coeff.	0.0508	-	-
## Cramer V	0.0508	0.0000	0.1076

The phi coefficient is a measure of the degree of association between two categorical variables and is interpretable as a correlation coefficient. It is derived from the χ^2 -statistic, but is free of the influence of the total sample size (Fleiss, 1981). Being independent of the sample size is a desirable quality because the χ^2 -statistic itself is sensitive to sample size. As the sample size increases, the χ^2 value will increase even if the cell proportions remain unchanged.

Pearson's contingency coefficient and Cramer's V are also derived from the chi-square and in the 2x2 table they are identical to the Phi coefficient (and similar to the Phi coefficient in interpretation). These three measures of degree of association are well suited for nominal variables in which the order of the levels is meaningless.

Cramer's V is useful for comparing multiple χ^2 test statistics and is generalizable across contingency tables of varying sizes. It is not affected by sample size and therefore is very useful in situations, where a statistically significant test result is suspected to be the result of a large sample size instead of any substantive relationship between the variables. It is interpreted as a measure of the relative strength of an association between two variables. It goes from 0 to 1, where 1 indicates strong association. In 2x2-tables the range is -1 to 1. The value of 0.0508 shows a very small, resp. no association between sex and hair colour at all.

The following are measures of ordinal association that consider whether the variable Y tends to increase as X increases: Gamma, Kendall's tau-b, Stuart's tau-c, and Somers' D. These measures are appropriate for ordinal variables, and they classify pairs of observations as concordant or discordant. A pair is concordant if the observation with the larger value of X also has the larger value of Y. A pair is discordant if the observation with the larger value of X has the smaller value of Y. Refer to Agresti (1996) and the other references cited in the discussion of each measure of association.

(We switch the example here, because our HairEyeColour variables aren't ordinal.)

```
(job <- matrix(c(16,19,9,8,17,11,14,60,56), nrow=3,
               dimnames=list("satisfaction"=c("high","medium","low"),
                             "security"=c("high","medium","low"))))

##           security
## satisfaction high medium low
##      high      16      8  14
##      medium    19     17  60
##      low       9     11  56

Desc(job, verbose="high")
##
##           estimate  lwr.ci  upr.ci
## ... (output skipped)
## Goodman Kruskal Gamma      0.3960 0.2103 0.5817
## Kendall Tau-b              0.2405 0.1206 0.3603
## Stuart Tau-c               0.2106 0.1038 0.3174
## Somers D C|R               0.2238 0.1123 0.3354
## Somers D R|C               0.2583 0.1242 0.3924
## Pearson Correlation        0.2742 0.1442 0.3950
## Spearman Correlation        0.2633 0.1327 0.3850
## ... (output skipped)
```

Gamma is recommended when there are lots of ties in the data. Tau-b is recommended for square tables.

The Pearson correlation coefficient and the Spearman rank correlation coefficient are also appropriate for ordinal variables. The Pearson correlation describes the strength of the linear association between the row and column variables, and it is computed using the row and column scores specified. The Spearman correlation is computed with rank scores.

The polychoric correlation is not reported, but can be calculated with the function `CorPolychor`. It also requires ordinal variables and assumes that the variables have an underlying bivariate normal distribution.

The measures of association lambda and uncertainty coefficient do not require ordinal variables, but they are appropriate for nominal variables.

Lambda has another concept than chi-squares. With Lambda the proportional reduction in error will be calculated. Lambda allows deciding, if the prediction of a class can be improved by using the other variable.

```

Desc(apply(Titanic, c(2,4), sum), verbose="high", rfrq="000")
## -----
... (output skipped)

##          estimate   lwr.ci   upr.ci
... (output skipped)
## Lambda C|R          0.3066  0.2568  0.3564
## Lambda R|C          0.0000  0.0000  0.0000
## Lambda sym          0.1846  0.1546  0.2146
## Uncertainty Coeff. C|R 0.1569  0.1283  0.1854
## Uncertainty Coeff. R|C 0.1903  0.1570  0.2237
## Uncertainty Coeff. sym 0.1720  0.1414  0.2026
## Mutual Information    0.1424    -    -
##
##      Survived
##      No   Yes   Sum
## Sex
## Male   1'364   367 1'731
## Female   126   344   470
## Sum     1'490   711 2'201

```

Without information about the sex, the best prediction for surviving would be “No”. We would guess 2201-1490=711 FALSE (Error E1=711) and 1490 correct. Using the variable sex we would guess survived “Yes” for women and “No” for men. So we would guess correct 344 women and 1364 men and 126 women and 367 men not correct (leading to an error E2=126+367=493). Lambda is then calculated as

$$\lambda(C|R) = \frac{E1 - E2}{E1} = \frac{711 - 493}{711} = 0.3066$$

The (C|R) notation indicates that the column variable is to be predicted by using the row variable. Thus, using the variable Sex (row-variable R) we make 30% less errors in predicting Survival of Titanic disaster (column variable C).

Note that we would not profit by the variable survived to predict the sex of a person, as the according lambda value R|C is 0.

Asymptotic confidence limits for all statistics are computed. The confidence coefficient is determined according to the value of the `conf.level` option, which by default equals 0.95 and produces 95% confidence limits.

13 Cases

The following cases are taken more or less verbatim from the SAS-Freq documentation^[4] and recalculated with base R and specific DescTools functions. The comments and descriptions have partly been adopted.

13.1 Eye colour - Binomial Proportions for One-Way Frequency Tables

The binomial proportions are computed as the proportion of observations for all the levels of the variable. The following statements compute the proportion of children with brown eyes (from the data set in Example 28.1 on page 1335) and test this value against the hypothesis that the proportion is 50%. Also, these statements test whether the proportion of children with fair hair is 28%.

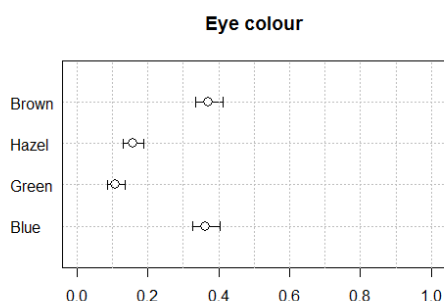
```
tab <- as.table(apply(HairEyeColor, 2, sum)[c("Brown", "Hazel", "Green", "Blue")])
Desc(tab)
```

```
## -----
## tab (table)
##
## Summary:
## n: 592, rows: 4
##
## Pearson's Chi-squared test (1-dim uniform):
## X-squared = 133.47, df = 3, p-value < 2.2e-16
##
##   level freq  perc cumfreq cumperc
## 1 Brown  220 37.2%    220    37.2%
## 2 Hazel   93 15.7%    313    52.9%
## 3 Green   64 10.8%    377    63.7%
## 4 Blue  215 36.3%    592   100.0%

xci <- BinomCI(tab, sum(tab))
rownames(xci) <- rownames(tab)
print(xci, digits=3)

##           est lwr.ci upr.ci
## Brown 0.372 0.3336 0.411
## Hazel 0.157 0.1300 0.189
## Green 0.108 0.0856 0.136
## Blue 0.363 0.3254 0.403
```

Let's produce a plot of that:



```
PlotDot(xci[,1], main="Eye colour", pch=NA,
        args.errbars = list(
          from=xci[,2], to=xci[,3],
          mid=xci[,1], pch=21, cex=1.4),
        xlim=c(0,1))

abline(v=seq(0,1,0.1), col="grey", lty="dotted")
```

PlotDot

Figure 13.1 Dotplot of marginal proportions for eye colour.

The estimation of simultaneously calculated confidence intervals for multinomial proportions according to the method of Sison and Glaz leads to slightly broader confidence intervals especially for the smaller groups (Hazel, Green).

```
print(MultinomCI(tab), digits=3)
```

```
##           est lwr.ci upr.ci
## Brown 0.372 0.3294 0.415
## Hazel 0.157 0.1149 0.201
## Green 0.108 0.0659 0.152
## Blue  0.363 0.3209 0.407
```

13.2 Cochran-Armitage Trend Test

In clinical trials, a dose response study is often conducted to investigate the relationship between increasing dosage and the effect of the drug under study. Usually the dose levels tested are ordinal, and the effect of the drug is measured in binary. In this case, Cochran-Armitage trend test is frequently used to test for trend among binomial proportions.

```
d.lungtumor <- data.frame(dose = rep(c(0, 1, 2), c(40, 50, 48)),
                          tumor = c(rep(c(0, 1), c(38, 2)),
                                    rep(c(0, 1), c(43, 7)),
                                    rep(c(0, 1), c(33, 15))))
lung <- table(d.lungtumor$dose, d.lungtumor$tumor)
Desc(lung, rfrq="010")
```

```
... (output skipped)
##           tumor
##           0    1    Sum
## dose
## 0   freq      38     2    40
##     p.row    95.0%  5.0%    .
## 1   freq      43     7    50
##     p.row    86.0% 14.0%    .
## 2   freq      33    15    48
##     p.row    68.8% 31.2%    .
## Sum freq      114    24   138
##     p.row      .     .     .
```

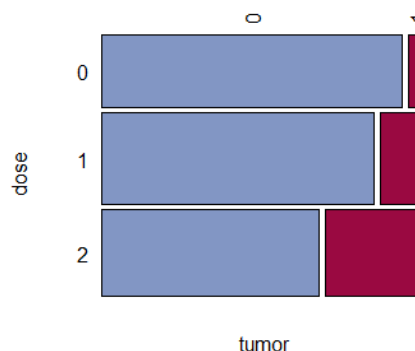


Figure 13.2 Lung cancer proportions.

```
CochranArmitageTest(lung, alternative = "increasing")
```

```
##           Cochran-Armitage test for trend
##
## data:  lung
## Z = -3.2735, dim = 3, p-value = 0.0005311
## alternative hypothesis: increasing
```

The Cochran-Armitage test supports the trend hypothesis. The small right-sided p-value (alternative = "increasing") indicate that the probability of the column 1 level (lungtumor = 1) increase as dose increases.

**Cochran
ArmitageTest**

13.3 Heart – 2x2-Table

This example computes chi-square tests and Fisher's exact test to compare the probability of coronary heart disease for two types of diet. It also estimates the relative risks and computes exact confidence limits for the odds ratio.

The data set contains hypothetical data for a case-control study of high fat diet and the risk of coronary heart disease. The data can be entered as:

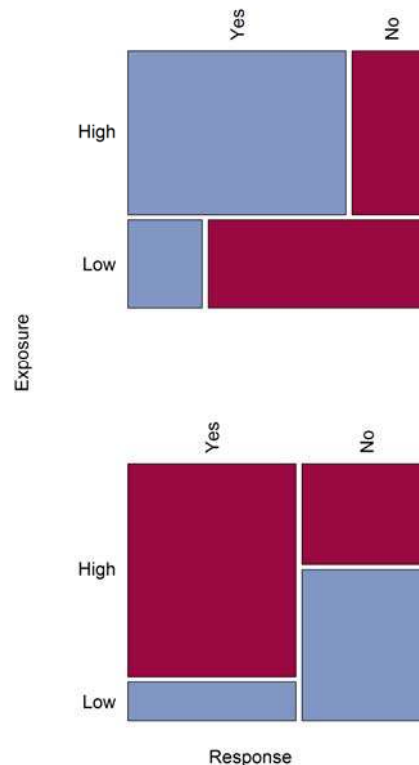
```
heart <- as.table(matrix(c(11, 2, 4, 6), nrow=2,
                        dimnames = list(Exposure = c("High", "Low"),
                                         Response = c("Yes", "No"))))
Label(heart) <- "Table of Response by Exposure"
```

The data is sorted in descending order by both variables, Exposure and Response, so that the first cell of the 2x2-table contains the frequency of positive exposure and positive response.

```
Desc(heart, main="Case-Control Study of High Fat/Cholesterol Diet")
```

will produce the following result:

```
## Case-Control Study of High Fat/Cholesterol Diet
##   Table of Response by Exposure
##
##
## Summary:
## n: 23, rows: 2, columns: 2
##
## Pearson's Chi-squared test (cont. adj):
## X-squared = 3.1879, df = 1, p-value = 0.0741
## Fisher's exact test p-value = 0.03931
## McNemar's chi-squared = 0.16667, df = 1,
## p-value = 0.6831
##
## Warning message:
## Exp. counts < 5: Chi-squared approx. may
## be incorrect!!
##
##               estimate lwr.ci upr.ci
##
## odds ratio          8.250  1.154 59.003
## rel. risk (col1)     2.933  0.850 10.120
## rel. risk (col2)     0.356  0.140  0.901
##
## Phi-Coefficient      0.464
## Contingency Coeff.    0.421
## Cramer's V           0.464
##
##           Response
##           Yes    No    Sum
## Exposure
## High   freq      11     4     15
##         perc     47.8%  17.4% 65.2%
##         p.row     73.3%  26.7% .
##         p.col     84.6%  40.0% .
##
## Low    freq       2     6     8
##         perc       8.7%  26.1% 34.8%
##         p.row     25.0%  75.0% .
##         p.col     15.4%  60.0% .
##
## Sum     freq      13     10     23
##         perc     56.5%  43.5% 100.0%
##         p.row     .     .     .
##         p.col     .     .     .
##
```



We learn that we have a total of 23 persons in our dataset and that the table has two rows and 2 columns. The association between the response and exposure appears not be existent, as the chi-square test is not significant ($p = 0.0741$).

However, if the expected value of one or more cells is less than 5, the chi-square test may not be valid. A specific warning indicates, that this is here the case. Fisher's exact test is an alternative test which does not depend on the expected values and is the appropriate test in this situation. It analyses whether the probability of heart disease in the high fat group differs from the one in the low fat group; since this p-value is small ($p < 0.05$), the alternative hypothesis is supported. Note that only the one-sided test will be reported.

The function expects the table to have the risk factor in rows and the response or outcome in the columns. The positive risk factor is preferred to be in the first row and the positive response in the first column:

Risk factor	Response	
	Yes	No
Yes	A	B
No	C	D

The odds ratio is then defined as

$$OR = \frac{\frac{A}{B}}{\frac{C}{D}} = \frac{A \cdot D}{B \cdot C} = \frac{11 \cdot 6}{4 \cdot 2} = 8.25$$

Recall that the odds of an event occurring is the ratio of p/q where p is the probability of the event occurring and q is the probability of the event not occurring. The odds ratio provides in fact an estimate of the relative risk when an event is rare (which here is not the case!).

The estimate indicates that the odds of heart disease are 8.25 times higher in the high fat diet group; however, the wide confidence limits (1.154, 59.003) indicate that this estimate has low precision.

The relative risk is the ratio of the probability of the heart disease occurring in the risk group (high fat diet) to the probability of the heart disease occurring in the comparison, non-exposed group (low fat diet). This is reported as rel. risk (col1) in the output above.

$$RR_1 = \frac{\frac{A}{A+B}}{\frac{C}{C+D}} = \frac{A \cdot (C+D)}{C \cdot (A+B)} = \frac{11 \cdot (2+6)}{2 \cdot (11+4)} = 2.933$$

A relative risk greater than 1 indicates that the probability of positive response is greater (here: heart disease) in row 1 (here: high fat diet group) than in row 2 (here: low fat diet group). Similarly, a relative risk less than 1 would indicate that the probability of positive response is less in row 1 than in row 2. The strength of association increases with the deviation from 1.

The relative risk column 2 uses the observations in this column to calculate the ratio.

$$RR_2 = \frac{\frac{B}{A+B}}{\frac{D}{C+D}} = \frac{B \cdot (C+D)}{D \cdot (A+B)} = \frac{4 \cdot (2+6)}{6 \cdot (11+4)} = 0.356$$

Recall an incidence rate is the proportion of new cases (outcomes) occurring over a period of any one time. Therefore the risk of an outcome makes sense in the context of prospective cohort studies where the outcome has not occurred in any case at the start of the study.

While the relative risk RR is a measure which is appropriate for prospective cohort studies, the odds ratio OR can be used for crosssectional case-control studies as well as prospective studies. In both cases, a value of 1 indicates no difference between groups.

Interchanging the row and column variables or modifying the table order will result in different values of odds ratio and relative risks. Reversing the columns for instance will result in the reciprocal OR:

```
OddsRatio(heart)
## [1] 8.25
```

```
1 / OddsRatio(Rev(heart, 1))  
## [1] 8.25
```

The interpretations should however remain consistent.

13.4 Optalmological Test – Confusion Matrix

A 2x2 table is often used in epidemiology when it comes to compare new diagnostic tests. Usually new tests need to be compared against a 'gold' standard. The gold standard stands for the best single test (or a combination of tests) that is considered the current preferred method of diagnosing a particular disease.

Let's take an example from (Rajul, 2008). One hundred persons with primary angle closure glaucoma are examined by a new test. Seventy-five of them had the disease.

Test	Disease	
	Yes	No
Yes	A (TP)	B (FP)
No	C (FN)	D (TN)

TP: True positive, FP: False positive, FN: False negative, TN: True negative

The two columns indicate the actual condition of the subjects, diseased or non-diseased. The rows indicate the results of the test, positive or negative. Cell A contains true positives, subjects with the disease and positive test results. Cell D subjects do not have the disease and the test agrees.

A good test will have minimal numbers in cells B and C. Cell B identifies individuals without disease but for whom the test indicates 'disease'. These are false positives. Cell C has the false negatives.

Sensitivity is the probability that a test will indicate 'disease' among those with the disease. Specificity is the fraction of those without disease who will have a negative test result. Positive and negative predictive values are influenced by the prevalence of disease in the population that is being tested. If we test in a high prevalence setting, it is more likely that persons who test positive truly have disease than if the test is performed in a population with low prevalence.

(Find the other statistics explained in the help file of the function Conf.)

```
opht <- as.table(matrix(c(75, 25, 15, 85), nrow=2,
                        dimnames = list(Test = c("Pos", "Neg"),
                                         Disease = c("Pos", "Neg"))))
```

```
Conf(opht)
```

Conf

```
## Confusion Matrix and Statistics
```

```
##
```

```
##           Reference
```

```
## Prediction Pos Neg
```

```
##           Pos  75  15
```

```
##           Neg  25  85
```

```
##
```

```
##           Accuracy : 0.8000
```

```
##           95% CI : (0.7391, 0.8495)
```

```
##           No Information Rate : 0.5000
```

```
##           P-Value [Acc > NIR] : < 2.2e-16
```

```
##
```

```
##           Kappa : 0.6000
```

```
## Mcnemar's Test P-Value : 0.1547
```

```
##
```

```
##           Sensitivity : 0.7500
```

```
##           Specificity : 0.8500
```

```
##           Pos Pred Value : 0.8333
```

```
##           Neg Pred Value : 0.7727
```

```
##           Prevalence : 0.5000
```

```
##           Detection Rate : 0.3750
```

```
##           Detection Prevalence : 0.5000
```

```
##           Balanced Accuracy : 0.8000
```

```
##           F-val Accuracy : 0.7895
```

```
##
```

```
##           'Positive' Class : Pos
```

Sensitivity:

$$\frac{A}{A+C} = \frac{75}{75+25} = 0.75$$

Specificity:

$$\frac{D}{B+D} = \frac{85}{85+15} = 0.85$$

Pos. Predicted Value:

$$\frac{A}{A+B} = \frac{75}{75+15} = 0.83$$

Neg. Predicted Value:

$$\frac{D}{D+C} = \frac{85}{85+25} = 0.77$$

13.5 Skin - Agreement Study

Medical researchers are interested in evaluating the efficacy of a new treatment for a skin condition. Dermatologists from participating clinics were trained to conduct the study and to evaluate the condition. After the training, two dermatologists examined patients with the skin condition from a pilot study and rated the same patients. The possible evaluations are terrible, poor, marginal, and clear.

In order to evaluate the agreement of the diagnoses (a possible contribution to measurement error in the study), the kappa coefficient is computed.

```
ParseSASDataLines("
  data d.SkinCondition;
    input Derm1 $ Derm2 $ Count;
  datalines;
  terrible terrible 10   terrible poor 4   terrible marginal 1   terrible clear 0
  poor terrible 5       poor poor 10      poor marginal 12      poor clear 2
  marginal terrible 2   marginal poor 4   marginal marginal 12  marginal clear 5
  clear terrible 0     clear poor 2     clear marginal 6     clear clear 13
;")
skin <- xtabs(Count ~ ., d.SkinCondition)
```

The function Agree computes raw simple percentage agreement among raters.

```
Agree(Untable(skin))
```

```
## [1] 0.5113636
```

```
## attr(,"subjects")
```

```
## [1] 88
```

```
## attr(,"raters")
```

```
## [1] 2
```


We learn that 51.1% of the ratings were the same between the two researchers. A less coarse approach to measure agreement is Cohen's kappa.

```
CohenKappa(skin, conf.level=0.95)
```

```
##      kappa      lwr.ci      upr.ci
## 0.3448753 0.2048513 0.4848994
```

```
CohenKappa(skin, conf.level=0.95, weights="Fleiss-Cohen")
```

```
##      kappa      lwr.ci      upr.ci
## 0.6607229 0.4207465 0.9006993
```

CohenKappa

The kappa coefficient has the value 0.3449, which indicates slight agreement between the dermatologists. The conclusion to reject the null hypothesis of no agreement is supported by the confidence interval of (0.2030, 0.4868), which suggests that the true kappa is greater than zero. The weighted kappa coefficient can be calculated by defining the `weights` argument. Its value is even larger (0.6607) than the unweighted kappa.

The Bowker's test for symmetry (reported by `mcnemar.test`) is not defined here (because of the zeros in the table).

13.6 Migraine - Statistics for a Stratified 2x2-Table

The data set `Migraine` contains hypothetical data for a clinical trial of migraine treatment. Subjects of both genders receive either a new drug therapy or a placebo. Their response to treatment is coded as 'Better' or 'Same'. The data are recorded as cell counts, and the number of subjects for each treatment and response combination is recorded in the variable `Count`. The following statements create a three-way table stratified by Gender, where Treatment forms the rows and Response forms the columns.

```
ParseSASDataLines("
  data d.Migraine;
    input Gender $ Treatment $ Response $ Count @@;
  datalines;
    female Active Better 16 female Active Same 11
    female Placebo Better 5 female Placebo Same 20
    male Active Better 12 male Active Same 16
    male Placebo Better 7 male Placebo Same 19
  ;
")
migraine <- xtabs(Count ~ Treatment + Response + Gender, d.Migraine)
```

How does this look like?

```
ftable(migraine, col.vars = c(1,3))
```

```
##      Treatment Active      Placebo
##      Gender  female male  female male
## Response
##   Better           16   12         5    7
##   Same            11   16        20   19
```

It's always a good idea to have a plot of the situation:

```
d.frm <- as.data.frame(prop.table(migraine, c(2,3)))
d.frm$Treatment <- reorder.factor(d.frm$Treatment, new.order =
c("Placebo","Active"))
d.frm$Response <- reorder.factor(d.frm$Response, new.order = c("Same","Better"))

library(lattice)
barchart(Freq ~ Response | Treatment + Gender, data=d.frm,
         col="steelblue",
         panel = function(x, ...) {
           panel.grid(h=-1, v=0)
           panel.barchart(x, ...)
         },
         par.settings = list(strip.background=list(col="lightgrey"),
                             layout.heights=list(strip=1.45)),
         par.strip.text = list(col="black"),
         layout=c(2,2), cex.axis=2, ylim=c(0,1), xlab="Response", ylab="Percent",
         scales=list(tck=c(0.8,0.8), col="black", x=list(cex=1), y=list(cex=1)),
         main="Migraine")
```

barchart

This code yields:

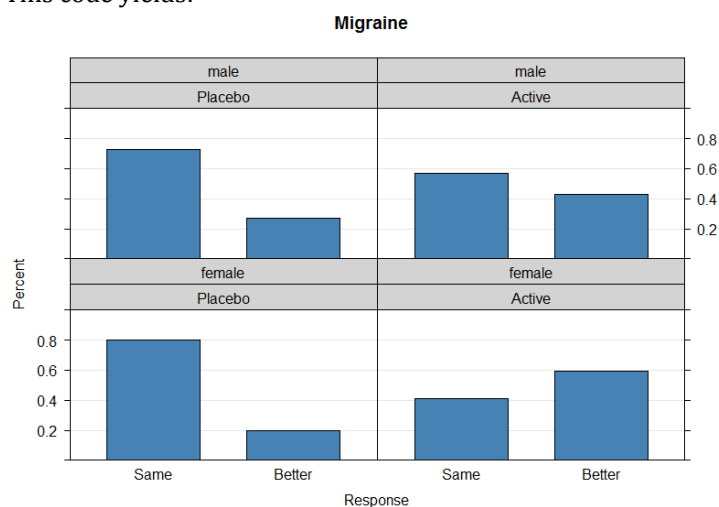


Figure 13.3 Trellis barplot of migraine patients.

The percentages are calculated so, that every panel has a total of 100%:

```
ptab <- prop.table(migraine, c(2,3))
ptab[] <- Format(ptab, digits=1, fmt="%")
ptab

## , , Treatment = Active

##      Gender
## Response female male
## Better 59.3% 42.9%
## Same   40.7% 57.1%

## , , Treatment = Placebo

##      Gender
## Response female male
## Better 20.0% 26.9%
## Same   80.0% 73.1%
```

Apparently the treatment seems to have an obvious effect. But the plot seems as well to indicate a gender effect, as the treatment is more pronounced for women than for men.

The function `mantelhaen.test` produces the Cochran-Mantel-Haenszel statistics. For this stratified 2x2 table, an estimate of the common odds ratio including its confidence interval is also displayed. (Note that the function expects the third dimension to be the strata, here gender.)

```
mantelhaen.test(migraine, alternative = "two.sided", correct = FALSE)

##      Mantel-Haenszel chi-squared test without continuity correction
##
## data:  migraine
## Mantel-Haenszel X-squared = 8.3052, df = 1, p-value = 0.003953
## alternative hypothesis: true common odds ratio is not equal to 1
## 95 percent confidence interval:
##  1.445613 7.593375
## sample estimates:
## common odds ratio
##      3.313168
```

mantelhaen.test

The significant p-value (0.004) indicates that the association between treatment and response remains strong after adjusting for gender.

A table of relative risks can be produced with

```
apply(migraine, 3, function(x) list(rbind(
  "Case-control (odds ratio)" = OddsRatio(x, conf.level = 0.95),
  "Cohort (col1 risk)"       = RelRisk(x, conf.level = 0.95),
  "Cohort (col2 risk)"       = RelRisk(Rev(x, 1), conf.level = 0.95))))

## $female
## $female[[1]]
##               odds ratio   lwr.ci   upr.ci
## Case-control (odds ratio)  5.818182 1.6755251 20.2033617
## Cohort (col1 risk)         2.962963 1.3713759 7.0036872
## Cohort (col2 risk)         0.337500 0.1427819 0.7291947
##
## $male
## $male[[1]]
##               odds ratio   lwr.ci   upr.ci
## Case-control (odds ratio)  2.0357143 0.6477707 6.397531
## Cohort (col1 risk)         1.5918367 0.7662184 3.454346
## Cohort (col2 risk)         0.6282051 0.2894904 1.305111
```

Because this is a prospective study, the relative risk estimate assesses the effectiveness of the new drug; the “Cohort (col1 risk)” values are the appropriate estimates for the first column, or the risk of improvement. The probability of migraine improvement with the new drug is just over two times the probability of improvement with the placebo.

The function `mantelhaen.test` displays also an estimate of the common odds ratio. This figure is calculated as [Agresti, p. 234]:

```
sum(apply(migraine, 3, function(x) prod(diag(x))/sum(x))) /
  sum(apply(migraine, 3, function(x) prod(diag(Rev(x, 1)))/sum(x)))

## 3.313168
```

The Breslow-Day test for homogeneity of the odds ratios can be calculated with the eponymous function. It tests the null hypothesis that the odds ratios for the q strata are all equal.

```
BreslowDayTest(migraine)

##      Breslow-Day Test on Homogeneity of the Odds Ratios
##
## data:  migraine
## X-squared = 1.4965, df = 1, p-value = 0.2212
```

BreslowDayTest

The large p-value (0.2212) indicates no significant gender difference in the odds ratios. Had the test for homogeneity of the odds ratios been statistically significant, a closer examination of each 2x2 table at each strata of the stratification variable would be required before making any further interpretations or conclusions.

Caution: Unlike the Cochran-Mantel-Haenszel statistics, the Breslow-Day test requires a large sample size within each stratum, and this limits its usefulness. In addition, the validity of the Cochran-Mantel-Haenszel tests does not depend on any assumption of homogeneity of the odds ratios; therefore, the Breslow-Day test should never be used as such an indicator of validity.

Homogeneity could also be assessed using Woolf's test.

```
WoolfTest(migraine)
```

WoolfTest

```
##      Woolf-test on Homogeneity of Odds Ratios (no 3-Way assoc.)
##
## data:  migraine
## X-squared = 1.4808, df = 1, p-value = 0.2236
```

Here the Woolf gives almost equivalent results to the BreslowDay test for consistency for the odds ratio.

The odds ratio for the treatment is

```
tab <- t(apply(migraine, c(1,2), sum))
OddsRatio(tab, conf.level = 0.95)
```

```
## odds ratio      lwr.ci      upr.ci
##    3.370370    1.461559    7.772108
```

Now, let's create logistic regression models on the raw data, first using just the two covariates Treatment and Gender:

```
r.glm <- glm(Response ~ Treatment + Gender, data=d.mig, family="binomial")
summary(r.glm)
```

glm

```
## Call:
## glm(formula = Response ~ Treatment + Gender, family = "binomial",
##      data = d.mig)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.2455  -1.0502  -0.6943   1.1108   1.7556
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)   -1.0602     0.3864  -2.744  0.00607 **
## TreatmentActive  1.2188     0.4271   2.853  0.00433 **
## Gendermale    -0.2398     0.4186  -0.573  0.56674
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 140.50  on 105  degrees of freedom
## Residual deviance: 131.55  on 103  degrees of freedom
## AIC: 137.55
##
## Number of Fisher Scoring iterations: 4
```

The estimates of the odds ratio are:

```
exp(coef(r.glm))

##      (Intercept) TreatmentActive  Gendermale
##      0.3463854      3.3830977      0.7867777
```

We learn that the treatment is significantly effective. Persons with treatment are 3.3 times as likely to report a positive response. The gender is not significant. By the way, also an interaction term would not become significant (not shown here).

14 References

- (1) Agresti A. (2002) Categorical Data Analysis. John Wiley & Sons.
- (2) Dalgaard P. (2008) Introductory Statistics with R (2. Aufl.), London, UK: Springer.
- (3) Friendly M. (2013) Working with categorical data with R and the vcd and vcdExtra packages, York University, Toronto.
<http://cran.r-project.org/web/packages/vcdExtra/vignettes/vcd-tutorial.pdf>
- (4) One-Way Frequency Tables using SAS PROC FREQ © TexaSoft, 2006
<http://www.stattutorials.com/SAS/TUTORIAL-PROC-FREQ-1.htm>
- (5) Presnell B. (2011) Course Notes sta4504-2011sp,
<http://www.stat.ufl.edu/~presnell/Courses/sta4504-2011sp/Notes/icda-notes-3x2.pdf>
- (6) SAS/STAT® 9.2 User's Guide, Second Edition, The FREQ Procedure (Book Excerpt) (2009)
<http://support.sas.com/documentation/cdl/en/statugfreq/63124/PDF/default/statugfreq.pdf>
- (7) Rajul Parikh et. al. Understanding and using sensitivity, specificity and predictive values, 2008 Indian J Ophthalmol. Jan-Feb; 56(1): 45–50.