

# hyperSpec Introduction

Claudia Beleites <[cbeleites@units.it](mailto:cbeleites@units.it)>  
CENMAT, DMRN, University of Trieste

March 3, 2011

## Reproducing the Examples in this Vignette

All spectra used in this manual are installed automatically with *hyperSpec*.  
Note that some definitions are executed in `vignette.defs`.

## Contents

<b>1. Introduction</b>	<b>3</b>
1.1. Notation . . . . .	3
<b>2. Remarks on R</b>	<b>4</b>
2.1. Generic Functions . . . . .	4
2.2. Functionality Can be Extended at Runtime . . . . .	4
2.3. Validity Checking . . . . .	4
2.4. Special Function Names . . . . .	4
2.4.1. The Names of Operators . . . . .	4
2.4.2. Assignment Functions . . . . .	5
<b>3. Loading and the package and configuration</b>	<b>5</b>
<b>4. The structure of hyperSpec objects</b>	<b>5</b>
<b>5. Functions provided by hyperSpec</b>	<b>6</b>
<b>6. Obtaining Basic Information about hyperSpec Objects</b>	<b>6</b>
<b>7. Creating a hyperSpec Object, Data Import and Export</b>	<b>7</b>
7.1. Creating a <i>hyperSpec</i> Object from Spectra Matrix and Wavelength Vector . . . . .	7
<b>8. The Logbook</b>	<b>8</b>
<b>9. Combining and Decomposing hyperspec Objects</b>	<b>8</b>
9.1. Binding Objects together . . . . .	8
9.2. Binding Objects that do not Share the Same Extra Data and/or Wavelength Axis . . . . .	9
9.3. Binding Objects that do not Share the Same Spectra . . . . .	9
9.4. Matrix Multiplication . . . . .	11
9.5. Decomposition . . . . .	11

<b>10. Access to the data</b>	<b>11</b>
10.1. Selecting and Deleting Spectra . . . . .	12
10.1.1. Random Samples . . . . .	12
10.1.2. Sequences . . . . .	13
10.2. Selecting Extra Data Columns . . . . .	13
10.3. Selecting Wavelength Ranges . . . . .	14
10.4. Deleting Wavelength Ranges . . . . .	15
10.4.1. Converting Wavelengths to Indices and vice versa . . . . .	15
10.4.2. Changing the Wavelength Axis . . . . .	16
10.4.3. Ordering the Wavelength Axis . . . . .	17
10.5. More on the Square-Bracket Operators for Replacing Values . . . . .	17
10.6. Fast Access to Parts of the <i>hyperSpec</i> Object . . . . .	18
10.7. Conversion to Long-Format <i>data.frame</i> . . . . .	18
<b>11. Plotting</b>	<b>19</b>
<b>12. Spectral (Pre)processing</b>	<b>19</b>
12.1. Cutting the Spectral Range . . . . .	19
12.2. Shifting Spectra . . . . .	19
12.2.1. Calculating the Shift . . . . .	20
12.3. Smoothing Interpolation . . . . .	21
12.4. Background Correction . . . . .	22
12.5. Offset Correction . . . . .	22
12.6. Baseline Correction . . . . .	22
12.7. Intensity Calibration . . . . .	23
12.7.1. Correcting by a constant, e. g. Readout Bias . . . . .	23
12.7.2. Correcting Wavelength Dependence . . . . .	23
12.7.3. Spectra Dependent Correction . . . . .	23
12.8. Normalization . . . . .	23
12.9. Centering the Data . . . . .	24
12.10 Variance Scaling . . . . .	24
12.11 Multiplicative Scatter Correction (MSC) . . . . .	25
12.12 Spectral Arithmetic . . . . .	25
<b>13. Data Analysis</b>	<b>25</b>
13.1. Data Analysis Methods using a <i>data.frame</i>	
e. g. Principal Component Analysis with <i>prcomp</i> . . . . .	25
13.1.1. PCA as Noise Filter . . . . .	26
13.2. Data Analysis using long-format <i>data.frame</i>	
e. g. plotting with <i>ggplot2</i> . . . . .	27
13.3. Data Analysis Methods using a matrix	
e. g. Hierarchical Cluster Analysis . . . . .	27
13.4. Calculating group-wise Sum Characteristics	
e. g. Cluster Mean Spectra . . . . .	28
13.5. Splitting an Object, and Binding a List of <i>hyperSpec</i> Objects . . . . .	28
<b>14. Speed Considerations</b>	<b>29</b>

**Suggested Packages***pls*: available*ggplot2*: available**1. Introduction**

*hyperSpec* is a R package that allows convenient handling of (hyper)spectral data sets, i. e. data sets comprising spectra together with further data on a per-spectrum basis. The spectra can be anything that is recorded over a common discretized axis.

*hyperSpec* works with any data that fits in this general scheme, so that the three terms may also be used for:

wavelength: frequency, wavenumbers, chemical shift, Raman shift,  $\frac{m}{z}$ , etc.

intensity: transmission, absorbance,  $\frac{e^-}{s}$ , ...

extra data: spatial information (spectral images, maps, or profiles), temporal information (kinetics, time series), concentrations (calibration series), class membership information, etc.  
Note that there is no restriction on the number of extra data columns.

Throughout the documentation of the package, the terms intensity and wavelength axes refer to the spectral ordinate and abscissa, respectively.

This vignette gives an introduction on basic working techniques using the R package *hyperSpec*.

*hyperSpec* comes with five data sets,

**chondro** a Raman map of chondrocytes in cartilage,

**flu** a set of fluorescence spectra of a calibration series, and

**laser** a time series of an unstable laser emission

**paracetamol** a Raman spectrum of paracetamol (acetaminophene) ranging from 100 to 3200  $\text{cm}^{-1}$  with overlapping wavelength ranges.

**barbituates** GC-MS spectra with differing wavelength axes as a list of 286 *hyperSpec* objects.

In this vignette, the data sets are used to illustrate appropriate procedures for different tasks and different spectra. In addition, the first three data sets are accompanied by their own vignettes showing exemplary work flows for the respective data type.

This document describes how to accomplish typical tasks in the analysis of spectra. It does not give a complete reference on particular functions. It is therefore recommended to look up the methods in R's help system using `? command`.

A list of all functions available in *hyperSpec* is given in appendix A (p. 30).

**1.1. Notation**

This vignette demonstrates working techniques mostly from a spectroscopic point of view: rather than going through the functions provided by *hyperSpec*, it is organized more closely on spectroscopic tasks. However, the functions discussed are printed on the margin for a fast overview.

In R, slots of a S4 class can be accessed directly by the `@` operator. In this vignette, the notation `@xxx` will thus mean “slot *xxx* of an object”. Likewise, named elements of a *list*, like the columns of a *data.frame*, are accessed by the `$` operator, and `$xxx` will be used for “column *xxx*”, and as an abbreviation for “column *xxx* of the *data.frame* in slot *data* of the object”.

## 2. Remarks on R

### 2.1. Generic Functions

*Generic Functions* are functions that apply to a wide range of data types or classes, e.g. *plot*, *print*, mathematical operators, etc. These functions can be implemented in a specialized way by each class. *hyperSpec* implements with a variety of such functions, see the table in appendix A on page 30.

### 2.2. Functionality Can be Extended at Runtime

R’s concept of functions offers much flexibility. Functions may be added or changed by the user in his *workspace* at any time. This is also true for methods belonging to a certain class. Neither restart of R nor reloading of the package or anything the like is needed. If the original function resides in a namespace (as it is the case for all functions in *hyperSpec*), the original function is not deleted. It is just masked by the user’s new function but stays accessible via the `::` operator.

This offers the opportunity of easily writing specialized functions that are adapted to specific tasks. *hyperSpec*’s vignettes use this to set up special versions of the lattice graphics functions that are already wrapped in `print` (see also [R FAQ: Why do lattice/trellis graphics not work?](#)) and allow the code in the code chunks of the vignettes to be exactly what one would type during an interactive R session. For the code, check the `vignettes.defs` file accompanying all *hyperSpec* vignettes.

### 2.3. Validity Checking

S4 classes have a mechanism to define and enforce that the data actually stored in the object is appropriate for this class. In other words, there is a mechanism of *validity checking*.

The functions provided by *hyperSpec* check the validity of *hyperSpec* objects at the beginning, and – if the validity could be broken by inappropriate arguments – also before leaving the function.

It is highly recommended to use validity checking also for user-defined functions. In addition, non-generic functions should first ensure that the argument actually is a *hyperSpec* object. The two tasks are accomplished by:

```
> chk.hy (object)
> validObject (object)
```

The first line checks whether `object` is a *hyperSpec* object, the second checks its validity. Both functions return `TRUE` if the checks succeed, otherwise they raise an error and stop.

### 2.4. Special Function Names

#### 2.4.1. The Names of Operators

Operators such as `+`, `*`, `%`, etc. are in fact functions in R. Thus they can be handed over as arguments to other functions (particularly to the vectorization functions `*apply`, `sweep`, etc.). In this case the name of the function must be quoted: ``*`` is the recommended style (although `"*"` will often work as well), e.g.:

slot	get	set
@wavelength	wl	wl<-
@data	[, [[, \$, as.data.frame, as.long.df, ...	[<-, [[<-, \$<-
@label	labels	labels<-
@log	logbook	logentry

Table 1: Get and set functions for the slots of *hyperSpec* objects

```
> sweep (flu, 2, mean, `~`)
```

These functions can also be called in a more function-like style (prefix notation):

```
> `+` (3, 5)
```

```
[1] 8
```

### 2.4.2. Assignment Functions

R allows the definition of functions that do an assignment (set some part of the object), such as:

```
> wl (flu) <- new.wavelength.values
```

an assignment to variable `wl`: `wl<-``.

## 3. Loading and the package and configuration

To load *hyperSpec*, use

```
> library (hyperSpec)
```

The global behaviour of *hyperSpec* can be configured via options. The values of the options are retrieved with `hy.getOptions` and `hygetOption`, and changed with `hy.setOptions`.

Currently, the only option provided is `log`, a logical specifying whether assignment functions should automatically add entries to the logbook (see section 8, p. 8).

## 4. The structure of *hyperSpec* objects

*hyperSpec* is a S4 (or new-style) class. Four slots contain the parts of the object:

`@wavelength` containing a numeric vector with the wavelength axis of the spectra.

`@data` a *data.frame* with the spectra and all further information belonging to the spectra

`@label` a list with appropriate labels (particularly for axis annotations)

`@log` a *data.frame* keeping track of what is done with the object

While the parts of the *hyperSpec* object can be accessed directly, it is good practice to use the functions provided by *hyperSpec* to handle the objects rather than accessing the slots directly (tab. 1). This also ensures that proper (*valid*) objects are retained.

Most of the data is stored in `@data`. This *data.frame* has one special column, `$spc`. It is the column that actually contains the spectra. The spectra are stored in a matrix inside this column, as illustrated in figure 1. Even if there are no spectra, `$spc` must still be present. It is then a matrix with zero columns.

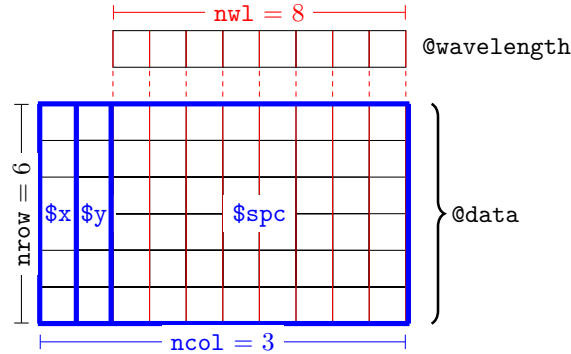


Figure 1: The structure of the data in a *hyperSpec* object.

Slot `@label` contains an element for each of the columns in `@data` plus one holding the label for the wavelength axis, `.wavelength`. The elements of the list may be anything suitable for axis annotations, i.e. they should be either character strings or expressions for “pretty” axis annotations (see e.g. figure 5 on page 22). To get familiar with expressions for axis annotation, see `? plotmath` and `demo (plotmath)`.

## 5. Functions provided by hyperSpec

Table A (p. 30) in the appendix gives an overview of the functions implemented by *hyperSpec*.

## 6. Obtaining Basic Information about hyperSpec Objects

As usual, the *print* and *show* methods display information about the object, and *summary* yields some additional details about the data handling done so far:

`print, show,`  
`summary`

```
> chondro

hyperSpec object
  875 spectra
  4 data columns
  300 data points / spectrum
wavelength: Delta * tilde(nu)/cm^-1 [numeric] 602 606 ... 1798
data: (875 rows x 4 columns)
  1. y: y/(mu * m) [numeric] -4.77 -4.77 ... 19.23
  2. x: x/(mu * m) [numeric] -11.55 -10.55 ... 22.45
  3. clusters: clusters [factor] matrix matrix ... lacuna + NA
  4. spc: I / a.u. [matrix300] 501.82 500.46 ... 169.29

> summary (chondro)

hyperSpec object
  875 spectra
  4 data columns
  300 data points / spectrum
wavelength: Delta * tilde(nu)/cm^-1 [numeric] 602 606 ... 1798
data: (875 rows x 4 columns)
  1. y: y/(mu * m) [numeric] -4.77 -4.77 ... 19.23
  2. x: x/(mu * m) [numeric] -11.55 -10.55 ... 22.45
  3. clusters: clusters [factor] matrix matrix ... lacuna + NA
  4. spc: I / a.u. [matrix300] 501.82 500.46 ... 169.29
log:
```

	short	long		date	user
1	.local	examp...	2011-03-03	15:31:03	cb@cb

The data set `chondro` consists of 875 spectra with 300 data points each, and 4 data columns: two for the spatial information, one factor with the results of a cluster analysis plus `$spc`. These informations can be directly obtained by

```
> nrow (chondro)
```

```
[1] 875
```

```
> nwl (chondro)
```

```
[1] 300
```

```
> ncol (chondro)
```

```
[1] 4
```

```
> dim (chondro)
```

```
nrow ncol nwl
875    4   300
```

The names of the columns in `@data` are accessed by

```
> colnames (chondro)
```

```
[1] "y"      "x"      "clusters" "spc"
```

```
colnames,
rownames,
dimnames, w1
```

Likewise, `rownames` returns the names assigned to the spectra, and `dimnames` yields a list of these three vectors (including also the column names of `$spc`). The column names of the spectra matrix are the wavelengths. They are accessed by `w1`, see section [10.4.2](#).

Extra data column names and rownames of the object may be set by `colnames<-` and `rownames<-`, respectively. `colnames<-` renames the labels as well.

```
colnames<-,
rownames<-
```

## 7. Creating a hyperSpec Object, Data Import and Export

*hyperSpec* comes with filters for a variety of file formats. These are discussed in detail in a separate vignette accessible via `vignette ("file-io")`.

### 7.1. Creating a hyperSpec Object from Spectra Matrix and Wavelength Vector

If the data is in R's workspace, a *hyperSpec* object is created by:

```
> spc <- new ("hyperSpec", spc = spectra.matrix, wavelength = wavelength.vector, data = extra.data)
```

The most frequently needed arguments are:

`spc` the spectra matrix

`wavelength` the wavelength axis vector

`data` the extra data (can already contain the spectra matrix in column `$spc`)

`label` a list with the proper labels. Do not forget the wavelength axis label in `$.wavelength` and the spectral intensity axis label in `$spc`.

## 8. The Logbook

Slot `@log` of *hyperSpec* objects is intended to keep track of the history of the object. This logbook part of the output of the `summary`, and can also be retrieved by `logbook`.

```
> logbook (flu)
```

	short.description	long.description	date	user
1	scan.txt.PerkinElmer	..., raw....	2011-02-13 21:50:53	cb@cb
2	\$<-	c, : [n....	2011-02-13 21:50:53	cb@cb
3	labels<-	c, c / (....	2011-02-13 21:50:53	cb@cb

New entries can be created manually by calling `logentry`:

logentry,  
logbook

```
> tmp <- logentry (flu, short = "test", long = "This could also be a list of parameters")
> logbook (tmp)
```

	short.description	long.description	date	user
1	scan.txt.PerkinElmer	..., raw....	2011-02-13 21:50:53	cb@cb
2	\$<-	c, : [n....	2011-02-13 21:50:53	cb@cb
3	labels<-	c, c / (....	2011-02-13 21:50:53	cb@cb
4	test	This cou....	2011-03-03 15:31:03	cb@cb

In addition, *hyperSpec* by default logs automatically all changes to the object:

```
> tmp <- tmp [1:3]
> logbook (tmp)
```

	short.description	long.description	date	user
1	scan.txt.PerkinElmer	..., raw....	2011-02-13 21:50:53	cb@cb
2	\$<-	c, : [n....	2011-02-13 21:50:53	cb@cb
3	labels<-	c, c / (....	2011-02-13 21:50:53	cb@cb
4	test	This cou....	2011-03-03 15:31:03	cb@cb
5	[]	i, j, .....	2011-03-03 15:31:03	cb@cb

The automatic logging mechanism can only log function calls and parameters (as opposed to the intention of the function call). *hyperSpec* functions that return a changed object allow to use more meaningful short descriptions: they are assigned via the argument *short*:

```
> tmp <- sweep (tmp, 2, mean, short = "centering")
> logbook (tmp)
```

	short.description	long.description	date	user
1	scan.txt.PerkinElmer	..., raw....	2011-02-13 21:50:53	cb@cb
2	\$<-	c, : [n....	2011-02-13 21:50:53	cb@cb
3	labels<-	c, c / (....	2011-02-13 21:50:53	cb@cb
4	test	This cou....	2011-03-03 15:31:03	cb@cb
5	[]	i, j, .....	2011-03-03 15:31:03	cb@cb
6	centering	MARGIN, ....	2011-03-03 15:31:03	cb@cb

## 9. Combining and Decomposing hyperspec Objects

### 9.1. Binding Objects together

*hyperspec* Objects can be bound together, either by columns (`cbind`) to append a new spectral range or by row (`rbind`) to append new spectra:

cbind rbind

```
> dim (flu)
```

```
nrow ncol nwl
6 3 181
```

```
> dim (cbind (flu, flu))
```



```
nrow ncol nwl
  6    3  362
```

```
> dim (rbind (flu, flu))
```

```
nrow ncol nwl
 12    3  181
```

There is also a more general function, `bind`, taking the direction ("`r`" or "`c`") as first argument followed by the objects to bind either in separate arguments or in a list.

As usual for `rbind` and `cbind`, the objects that should be bound together must have the same rows and columns, respectively.

## 9.2. Binding Objects that do not Share the Same Extra Data and/or Wavelength Axis

`collapse` combines objects that should be bound together by row, but they do not share the columns and/or spectral range. The resulting object has all columns from all input objects, and all wavelengths from the input objects. If an input object does not have a particular column or wavelength, its value in the resulting object is `NA`.

`collapse`

The `barbituates` data is a list of 286 *hyperSpec* objects, each containing one mass spectrum. The spectra have between 4 and 101 data points each.

```
> barb <- collapse (barbituates)
> wl (barb) [1 : 25]
```

```
[1] 160.90 158.85 147.00 140.90 133.05 130.90 119.95 119.15 118.05 116.95 112.90 106.00 105.10
[14] 98.95 96.95 91.00 85.05 83.05 77.00 71.90 71.10 70.00 69.00 57.10 56.10
```

The resulting object does not have an ordered wavelength axis. This can be obtained in a second step:

```
> barb <- orderwl (barb)
> barb [[1:3, , min ~ min + 10i]]
```

```
      25.95 26.05 26.15 26.95 27.05 27.15 28.05 28.15 29.05 29.15 29.95
[1,]    NA    NA    NA    NA   562    NA    NA 11511  6146    NA    NA
[2,]    NA    NA    NA    NA    NA   618 10151    NA  5040    NA    NA
[3,]    NA    NA    NA    NA   638    NA    NA 10722  5253    NA    NA
```

## 9.3. Binding Objects that do not Share the Same Spectra

`merge` adds a new spectral range (like `cbind`), but works also if spectra are missing in one of the objects. The arguments `by`, `by.x`, and `by.y` specify which columns should be used to decide which spectra are the same. The arguments `all`, `all.x`, and `all.y` determine whether spectra should be kept for the result set if they appear in only one of the objects. For details, see also the help on the base function `merge`.

`merge`

As an example, let's construct a version of the `chondro` data like being taken as two maps with different spectral ranges. In each data set, some spectra are missing.

```
> chondro.low <- sample (chondro [, , 600 ~ 1200], 700)
> nrow (chondro.low)

[1] 700

> chondro.high <- sample (chondro [, , 1400 ~ 1800], 700)
> nrow (chondro.high)
```

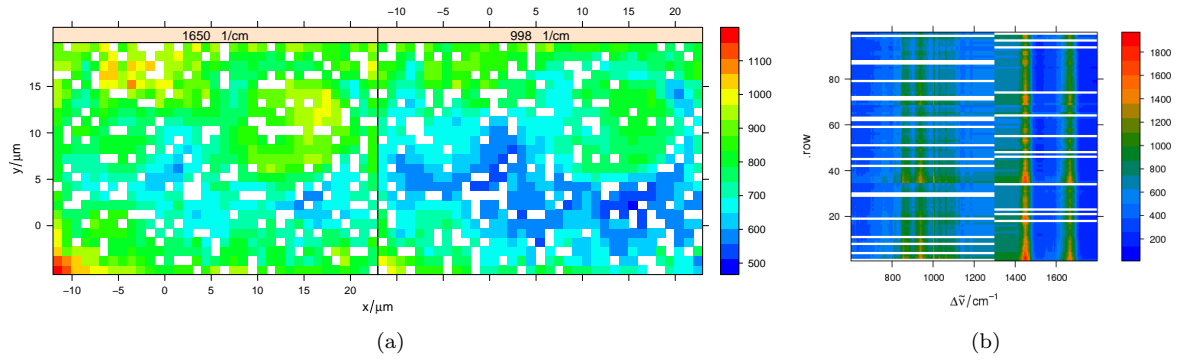


Figure 2: (a) For both spectral ranges some spectra are missing. (b) The missing parts of the spectra are filled with NA.

```
[1] 700
```

As all extra data columns are the same, no special declarations are needed for merging the data:

```
> chondro.merged <- merge (chondro.low, chondro.high)
> nrow (chondro.merged)
```

```
[1] 561
```

By default, the result consists of only those spectra, where *both* spectral ranges were available. To keep all spectra replacing missing parts by NA:

```
> chondro.merged <- merge (chondro.low, chondro.high, all = TRUE)
> nrow (chondro.merged)
```

```
[1] 839
```

```
> merged <- merge (chondro [1:7,, 610 ~ 620], chondro [5:10,, 615 ~ 625], all = TRUE)
> merged$.
```

	y	x	clusters	.nx	.ny	spc.610	spc.614	spc.618	spc.614	spc.618	spc.622	spc.626
1	-4.77	-11.55	matrix	1	NA	488.63	466.18	492.00	NA	NA	NA	NA
2	-4.77	-10.55	matrix	2	NA	489.48	465.05	490.53	NA	NA	NA	NA
3	-4.77	-9.55	matrix	3	NA	456.03	436.62	458.06	NA	NA	NA	NA
4	-4.77	-8.55	matrix	4	NA	464.82	444.85	470.02	NA	NA	NA	NA
5	-4.77	-7.55	matrix	5	1	428.66	410.80	433.12	410.80	433.12	461.19	397.38
6	-4.77	-6.55	matrix	6	2	426.07	407.86	431.21	407.86	431.21	458.15	394.18
7	-4.77	-5.55	lacuna	7	3	412.37	396.50	421.27	396.50	421.27	445.54	382.72
8	-4.77	-4.55	lacuna	NA	4	NA	NA	NA	381.95	406.25	429.67	368.46
9	-4.77	-3.55	lacuna	NA	5	NA	NA	NA	397.51	423.30	446.15	381.87
10	-4.77	-2.55	lacuna	NA	6	NA	NA	NA	377.39	402.23	424.19	362.43

If the spectra overlap, the result will have both data points. In the example here one could easily delete duplicate wavelengths. For real data, however, the duplicated wavelength will hardly ever contain the same values. The appropriate method to deal with this situation depends on the data at hand, but it will usually be some kind of spectral interpolation.

One possibility is removing duplicated wavelengths by using the mean intensity. This can conveniently be done by using `approx` using `method = "constant"`. For duplicated wavelengths, the intensities will be combined by the `tie` function. This already defaults to the mean, but we need `na.rm = TRUE`.

Thus, the function to calculate the new spectral intensities is

```
> approxfun <- function (y, wl, new.wl){
+   approx (wl, y, new.wl, method = "constant",
+         ties = function (x) mean (x, na.rm = TRUE)
+       )$y
+ }
```

which can be applied to the spectra:

```
> merged <- apply (merged, 1, approxfun,
+                 wl = wl (merged), new.wl = unique (wl (merged)),
+                 new.wavelength = "new.wl")
> merged$.
```

	y	x	clusters	.nx	.ny	spc.610	spc.614	spc.618	spc.622	spc.626
1	-4.77	-11.55	matrix	1	NA	488.6323....	466.1774....	492.0015....	NA	NA
2	-4.77	-10.55	matrix	2	NA	489.4758....	465.0506....	490.5328....	NA	NA
3	-4.77	-9.55	matrix	3	NA	456.0323....	436.6220....	458.0576....	NA	NA
4	-4.77	-8.55	matrix	4	NA	464.8207....	444.8485....	470.0171....	NA	NA
5	-4.77	-7.55	matrix	5	1	428.6619....	410.7955....	433.1227....	461.1903....	397.3773....
6	-4.77	-6.55	matrix	6	2	426.0734....	407.8569....	431.2144....	458.1502....	394.1775....
7	-4.77	-5.55	lacuna	7	3	412.3674....	396.5000....	421.2737....	445.5431....	382.7197....
8	-4.77	-4.55	lacuna	NA	4		NA 381.9504....	406.2470....	429.6728....	368.4599....
9	-4.77	-3.55	lacuna	NA	5		NA 397.5075....	423.3002....	446.1478....	381.8674....
10	-4.77	-2.55	lacuna	NA	6		NA 377.3917....	402.2348....	424.1901....	362.4296....

## 9.4. Matrix Multiplication

Two *hyperSpec* objects can be matrix multiplied by `%*%`. For an example, see the principal component analysis below (section 13.1 on page 25).

## 9.5. Decomposition

Matrix decompositions are common operations during chemometric data analysis. The results, e.g. of a principal component analysis are two matrices, the so-called scores and loadings. The results can have either the same number of rows as the spectra matrix they were calculated from (scores-like), or they have as many wavelengths as the spectra (loadings-like).

Both types of result objects can be “re-imported” into *hyperSpec* objects with function `decomposition`. A scores-like object retains all per-spectrum information (i.e. the extra data) while the spectra matrix and wavelength vector are replaced. A loadings-like object retains the wavelength information, while extra data is deleted (set to NA) unless the value is constant for all spectra.

A demonstration can be found in the principal component analysis example (section 13.1) on page 25.

## 10. Access to the data

The main functions to retrieve the data of a *hyperSpec* object are `[]` and `[][]`.

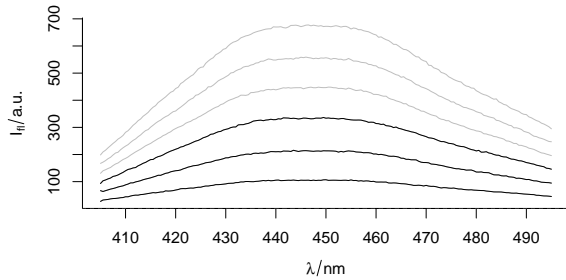
The difference between these functions is that `[]` returns a *hyperSpec* object, whereas the result of `[][]` is a *data.frame* if extra data columns were selected or otherwise the spectra matrix. Single extra data columns may be retrieved by `$`.

In order to change data, use `[]<-`, `[][]<-`, and `$<-` (see ).

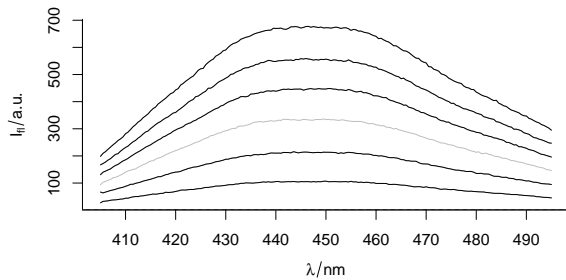
## 10.1. Selecting and Deleting Spectra

The extraction function `[]` takes the spectra as first argument (For detailed help: see `?`[]``). It may be a vector giving the indices of the spectra to extract (select), a vector with negative indices indicating which spectra should be deleted, or a logical. Note that a matrix given to `[]` will be treated as a vector.

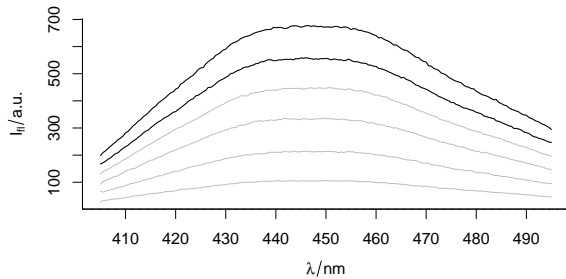
```
> plot (flu, col = "gray")
> plot (flu [1 : 3], add = TRUE)
```



```
> plot (flu, col = "gray")
> plot (flu [-3], add = TRUE)
```



```
> plot (flu, col = "gray")
> plot (flu [flu$c > 0.2], add = TRUE)
```



### 10.1.1. Random Samples

A random subset of spectra is conveniently selected by `sample` :

`sample`

```
> sample (chondro, 3)
```

```
hyperSpec object
3 spectra
4 data columns
```

```

      300 data points / spectrum
wavelength: Delta * tilde(nu)/cm^-1 [numeric] 602 606 ... 1798
data: (3 rows x 4 columns)
  1. y: y/(mu * m) [numeric] 1.23 3.23 19.23
  2. x: x/(mu * m) [numeric] 8.45 -6.55 11.45
  3. clusters: clusters [factor] matrix lacuna matrix
  4. spc: I / a.u. [matrix300] 292.90 291.55 ... 164.32

```

If appropriate indices into the spectra are needed instead, use `isample`:

`isample`

```

> isample (chondro, 3)

[1] 653 546 26

```

### 10.1.2. Sequences

Sequences of every  $n^{\text{th}}$  spectrum or the like can be retrieved with `seq`:

`seq`

```

> seq (chondro, length.out = 3, index = TRUE)

[1] 1 438 875

> seq (chondro, by = 100)

hyperSpec object
  9 spectra
  4 data columns
  300 data points / spectrum
wavelength: Delta * tilde(nu)/cm^-1 [numeric] 602 606 ... 1798
data: (9 rows x 4 columns)
  1. y: y/(mu * m) [numeric] -4.77 -2.77 ... 17.23
  2. x: x/(mu * m) [numeric] -11.55 18.45 ... 18.45
  3. clusters: clusters [factor] matrix matrix ... lacuna
  4. spc: I / a.u. [matrix300] 501.82 400.94 ... 124.64

```

Here, indices may be requested using `index = TRUE`.

## 10.2. Selecting Extra Data Columns

The second argument of the extraction functions `[]` and `[[[]]` specifies the (extra) data columns. They can be given like any column specification for a *data.frame*, i. e. numeric, logical, or by a vector of the column names:

They can be given like any column specification for a *data.frame*, i. e. numeric, logical, or by a vector of the column names:

```

> colnames (chondro)

[1] "y"      "x"      "clusters" "spc"

> chondro [[1 : 3, 1]]

      y
1 -4.77
2 -4.77
3 -4.77

> chondro [[1 : 3, -4]]

      y      x clusters
1 -4.77 -11.55  matrix
2 -4.77 -10.55  matrix
3 -4.77  -9.55  matrix

```

```
> chondro [[1 : 3, "x"]]

      x
1 -11.55
2 -10.55
3  -9.55

> chondro [[1 : 3, c (TRUE, FALSE)]]      # note the recycling!

      y clusters
1 -4.77  matrix
2 -4.77  matrix
3 -4.77  matrix
```

To select one column, the \$ operator is more convenient:

```
> flu$c

[1] 0.05 0.10 0.15 0.20 0.25 0.30
```

*hyperSpec* supports command line completion for the \$ operator.

The extra data may also be set this way:

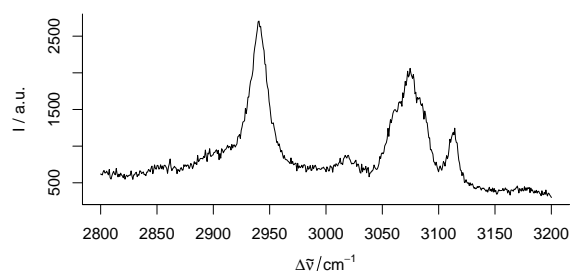
```
> flu$n <- list (1 : 6, label = "sample no.")
```

This function will append new columns, if necessary.

### 10.3. Selecting Wavelength Ranges

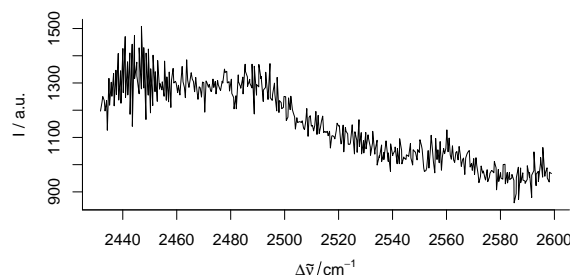
Wavelength ranges can easily be selected using []'s third argument:

```
> plot (paracetamol [, , 2800 ~ 3200])
```



By default, the values given are treated as wavelengths, if they are indices into the columns of the spectra matrix, use `wl.index = TRUE`:

```
> plot (paracetamol [, , 2800 : 3200, wl.index = TRUE])
```

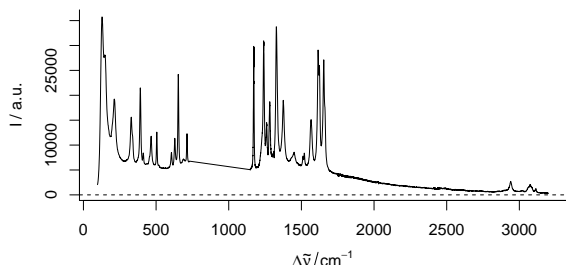


Section 10.4.1 (p. 15) details into the different possibilities of specifying wavelengths.

## 10.4. Deleting Wavelength Ranges

Deleting wavelength ranges may be accomplished using negative index vectors together with `wl.index = TRUE`.

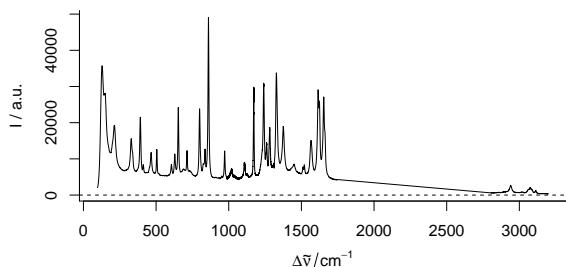
```
> plot (paracetamol [, -(500 : 1000)], wl.index = TRUE)
```



However, this mechanism works only if the proper indices are known.

If the range to be cut out is rather known in the units of the wavelength axis, it is easier to select the remainder of the spectrum instead. To delete the spectral range from 1750 to 2800  $\text{cm}^{-1}$  of the paracetamol spectrum one can thus use:

```
> plot (paracetamol [, c (min ~ 1750, 2800 ~ max)])
```



(It is possible to produce a plot of this data where the cut range is not bridged by a line and the wavelength axis is cut in order to save space. For details see the “plotting” vignette).

### 10.4.1. Converting Wavelengths to Indices and vice versa

Spectra in *hyperSpec* have always discretized wavelength axes, they are stored in a matrix with column corresponding to one wavelength. *hyperSpec* provides two conversion functions: `i2wl` returns the wavelength corresponding to the given indices and `wl2i` calculates index vectors from wavelengths.

wl2i i2wl

If the wavelengths are given as a numeric vector, they are each converted to the corresponding wavelength. In addition there is a more sophisticated possibility of specifying wavelength ranges using a *formula*. The basic syntax is *start ~ end*. This yields a vector *index of start : index of end*.

The result of the formula conversion differs from the numeric vector conversion in three ways:

- The colon operator for constructing vectors accepts only integer numbers, the tilde (for formulas) does not have this restriction.
- If the vector does not take into account the spectral resolution, one may get only every  $n^{\text{th}}$  point or repetitions of the same index:

```
> wl2i (flu, 405 : 410)
```

```
[1] 1 3 5 7 9 11

> wl2i (flu, 405 ~ 410)

[1] 1 2 3 4 5 6 7 8 9 10 11

> wl2i (chondro, 1000 : 1010)

[1] 100 101 101 101 101 102 102 102 102 103 103

> wl2i (chondro, 1000 ~ 1010)

[1] 100 101 102 103
```

- If the object's wavelength axis is not ordered, the formula approach will give weird results. In that (probably rare) case, use `orderwl` first to obtain an object with ordered wavelength axis.

`start` and `end` may contain the special variables `min` and `max` that correspond to the lowest and highest wavelengths of the object:

```
> wl2i (flu, min ~ 410)

[1] 1 2 3 4 5 6 7 8 9 10 11
```

Often, specifications like *wavelength  $\pm n$  data points* are needed. They can be given using complex numbers in the formula. The imaginary part is added to the index calculated from the wavelength in the real part:

```
> wl2i (flu, 450 - 2i ~ 450 + 2i)

[1] 89 90 91 92 93

> wl2i (flu, max - 2i ~ max)

[1] 179 180 181
```

To specify several wavelength ranges, use a list containing the formulas and vectors<sup>1</sup>:

```
> wl2i (flu, 450 - 2i ~ 450 + 2i)

[1] 89 90 91 92 93

> wl2i (flu, c (min ~ 406.5, max - 2i ~ max))

[1] 1 2 3 4 179 180 181
```

This mechanism also works for the wavelength arguments of `[]`, `[[ ]]`, and `plotspc`.

#### 10.4.2. Changing the Wavelength Axis

Sometimes wavelength axes need to be transformed, e. g. converting from wavelengths to frequencies. In this case, retrieve the wavelength axis vector with `wl`, convert each value of the resulting vector and assign the result with `wl<-`. Also the label of the wavelength axis may need to be adjusted.

`wl, wl<-`

As an example, convert the wavelength axis of `laser` to frequencies. As the wavelengths are in nanometers, and the frequencies are easiest expressed in terahertz, an additional conversion factor of 1000 is needed:

```
> laser
```

---

<sup>1</sup>Formulas are combined to a list by `c`.



```
hyperSpec object
  84 spectra
  2 data columns
  36 data points / spectrum
wavelength: lambda/nm [numeric] 404.58 404.62 ... 405.82
data: (84 rows x 2 columns)
  1. t: t / s [numeric] 0 2 ... 5722
  2. spc: I / a.u. [matrix36] 164.65 179.72 ... 112.09
```

```
> wavelengths <- wl (laser)
> frequencies <- 2.998e8 / wavelengths / 1000
> wl (laser) <- frequencies
> labels (laser, ".wavelength") <- "f / THz"
> laser
```

```
hyperSpec object
  84 spectra
  2 data columns
  36 data points / spectrum
wavelength: f / THz [numeric] 741.01 740.95 ... 738.76
data: (84 rows x 2 columns)
  1. t: t / s [numeric] 0 2 ... 5722
  2. spc: I / a.u. [matrix36] 164.65 179.72 ... 112.09
```

```
> rm (laser)
```

There are other possibilities of invoking `wl<-` including the new label, e.g.

```
> wl (laser, "f / THz") <- frequencies
```

and

```
> wl (laser) <- list (wl = frequencies, label = "f / THz")
```

see `?`wl<-`` for more information.

### 10.4.3. Ordering the Wavelength Axis

If the wavelength axis of an object needs reordering (e.g. after `collapse`), `orderwl` can be used:

`orderwl`

```
> barb <- collapse (barbituates [1 : 3])
> wl (barb)
```

```
[1] 160.90 158.85 147.00 140.90 133.05 130.90 119.95 119.15 118.05 116.95 112.90 106.00 105.10
[14] 98.95 96.95 91.00 85.05 83.05 77.00 71.90 71.10 70.00 69.00 57.10 56.10 55.00
[27] 43.85 43.05 41.10 40.10 39.00 32.15 31.15 30.05 29.05 28.15 27.05 132.95 131.00
[40] 120.05 119.05 117.95 113.00 105.90 82.95 72.00 69.10 56.00 44.05 40.00 30.15 28.05
[53] 27.15 84.15 68.90 55.10 43.95
```

```
> barb <- orderwl (barb)
> wl (barb)
```

```
[1] 27.05 27.15 28.05 28.15 29.05 30.05 30.15 31.15 32.15 39.00 40.00 40.10 41.10
[14] 43.05 43.85 43.95 44.05 55.00 55.10 56.00 56.10 57.10 68.90 69.00 69.10 70.00
[27] 71.10 71.90 72.00 77.00 82.95 83.05 84.15 85.05 91.00 96.95 98.95 105.10 105.90
[40] 106.00 112.90 113.00 116.95 117.95 118.05 119.05 119.15 119.95 120.05 130.90 131.00 132.95
[53] 133.05 140.90 147.00 158.85 160.90
```

## 10.5. More on the Square-Bracket Operators for Replacing Values

`[[ ]]` also accepts index matrices of size  $n \times 2$ . In this case, a vector of values from the spectra matrix is returned.

```

> indexmatrix <- matrix (c (1 : 3, 1 : 3), ncol = 2)
> indexmatrix

      [,1] [,2]
[1,]    1    1
[2,]    2    2
[3,]    3    3

> chondro [[indexmatrix, wl.index = TRUE]]

[1] 501.82 507.81 456.03

> diag (chondro [[1 : 3, , min ~ min + 2i]])

[1] 501.82 507.81 456.03

[[ ]] <- also accepts index matrices of size  $n \times 2$ .

> indexmatrix <- matrix (c (1 : 3, 1 : 3), ncol = 2)
> indexmatrix

      [,1] [,2]
[1,]    1    1
[2,]    2    2
[3,]    3    3

> chondro [[indexmatrix, wl.index = TRUE]]

[1] 501.82 507.81 456.03

> diag (chondro [[1 : 3, , min ~ min + 2i]])

[1] 501.82 507.81 456.03

```

## 10.6. Fast Access to Parts of the hyperSpec Object

[[ ]] \$. \$..

*hyperSpec* comes with three abbreviation functions for easy access to the data:

```

x [[ ]] returns the spectra matrix (x$spc).

x [[i, , l]] the cut spectra matrix is returned if wavelengths are specified in l.

x [[i, j, l]] If data columns are selected (second index), the result is a data.frame.

x [[i, , l]] <- Also, parts of the spectra matrix can be set (only indices for spectra and wave-
length are allowed for this function).

x [i, j] <- sets parts of x@data.

x $. returns the complete data.frame x@data, with the spectra in column $spc.

x $.. returns the extra data (x@data without x$spc).

x $.. <- sets the extra data (x@data without x$spc). However, the columns must match exactly
in this case.

```

## 10.7. Conversion to Long-Format data.frame

Some functions need the data being an *unstacked* or *long-format data.frame*. `as.long.df` is the appropriate conversion function. `as.long.df`

## 11. Plotting

*hyperSpec* offers a variety of possibilities to plot spectra, spectral maps, the spectra matrix, time series, depth profiles, etc.. This all is discussed in a separate document: see `vignette ("plotting")`.

## 12. Spectral (Pre)processing

### 12.1. Cutting the Spectral Range

□ □□

The extraction functions `[]` and `[] []` can be used to cut the spectra: Their third argument takes wavelength specifications as discussed above and also logicals (i.e. vectors specifying with TRUE/FALSE for each column of `$spc` whether it should be included or not).

`[]` returns a *hyperSpec* object, `[] []` the spectra matrix `$spc` (or the *data.frame* `@data` if in addition data columns were specified) only.

```
> flu[, , min ~ 408.5]
hyperSpec object
  6 spectra
  4 data columns
  8 data points / spectrum
wavelength: lambda/nm [numeric] 405.0 405.5 ... 408.5
data: (6 rows x 4 columns)
  1. file: [factor] rawdata/flu1.txt rawdata/flu2.txt ... rawdata/flu6.txt
  2. spc: I[fl]/a.u. [matrix8] 27.150 66.801 ... 256.89
  3. c: c / (mg / l) [numeric] 0.05 0.10 ... 0.3
  4. n: sample no. [integer] 1 2 ... 6

> flu[, , c (min ~ min + 2i, max - 2i ~ max)]
      405    405.5    406    494    494.5    495
[1,] 27.150 32.345 33.379 47.163 46.412 45.256
[2,] 66.801 63.715 66.712 96.602 96.206 94.610
[3,] 93.144 103.068 106.194 149.539 148.527 145.793
[4,] 130.664 139.998 143.798 201.484 198.867 195.867
[5,] 167.267 171.898 177.471 252.066 248.067 246.952
[6,] 198.430 209.458 215.785 307.519 302.325 294.649
```

### 12.2. Shifting Spectra

Sometimes, spectra need to be aligned along the spectral axis.

In general, two options are available for shifting spectra along the wavelength axis.

1. The wavelength axis can be shifted, while the intensities stay unaffected.
2. the spectra are interpolated onto a new wavelength axis, while the nominal wavelengths stay.

The first method is very straightforward:

```
> tmp <- chondro
> wl (tmp) <- wl (tmp) - 10
```

but it cannot be used if each spectrum (or groups of spectra) are shifted individually.

In that case, interpolation is needed. R offers many possibilities to interpolate (e.g. `approx` for constant / linear approximation, `spline` for spline interpolation, `loess` can be used to obtain smoothed approximations, etc.). The appropriate interpolation strategy will depend on the spectra, and *hyperSpec* therefore leaves it up to the user to select a sensible interpolation function.

As an example, we will use natural splines to do the interpolation. It is convenient to set it up as a function:

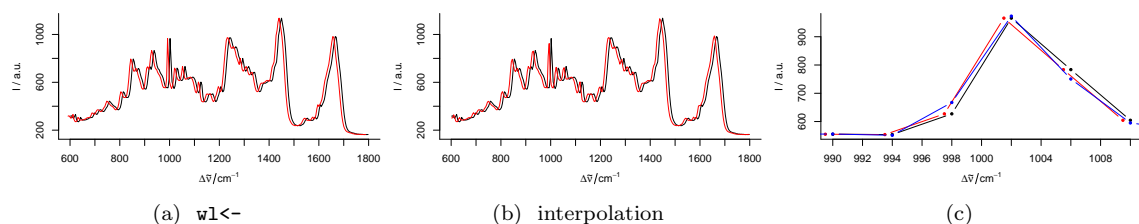


Figure 3: Shifting the Spectra along the Wavelength Axis. (a) Changing the wavelength values. (b) Interpolation. (c) Detail view of the phenylalanine band: shifting by `wl<-` (red) does not affect the intensities, while the spectrum is slightly changed by interpolations (blue).

```
> interpolate <- function (spc, shift, wl){
+   spline (wl + shift, spc, xout = wl, method = "natural")$y
+ }
```

This function can now be applied to a set of spectra:

```
> tmp <- apply (chondro, 1, interpolate, shift = -10, wl = wl (chondro))
```

If different spectra need to be offset by different shift, use a loop<sup>2</sup>

```
> shifts <- rnorm (nrow (chondro))
> tmp <- chondro [[]]
> for (i in seq_len (nrow (chondro)))
+   tmp [i, ] <- interpolate (tmp [i, ], shifts [i], wl = wl (chondro))
> chondro [[]] <- tmp
```

### 12.2.1. Calculating the Shift

Often, the shift in the spectra is determined by aligning a particular signal. This strategy works best with spectrally oversampled data that allows accurate determination of the signal position.

For the `chondro` data, let's use the maximum of the phenylalanine band between 990 and 1020  $\text{cm}^{-1}$ . As just the very maximum is too coarse, we'll use the maximum of a square polynomial fitted to the maximum and its two neighbours.

```
> find.max <- function (y, x){
+   pos <- which.max (y) + (-1:1)
+   X <- x [pos] - x [pos [2]]
+   Y <- y [pos] - y [pos [2]]
+   X <- cbind (1, X, X^2)
+   coef <- qr.solve (X, Y)
+   - coef [2] / coef [3] / 2 + x [pos [2]]
+ }
> bandpos <- apply (chondro [, 990 ~ 1020], 1, find.max, wl (chondro [, 990 ~ 1020]))
> refpos <- find.max (colMeans (chondro [, 990 ~ 1020]), wl (chondro [, 990 ~ 1020]))
> shift1 <- refpos - bandpos
```

A second possibility is to optimize the shift. For this strategy, the spectra must be sufficiently similar, while low spectral resolution is compensated by using larger spectral windows.

<sup>2</sup>`sweep` cannot be used here, and while there is the possibility to use `sapply` or `mapply`, they are not faster than the for loop in this case. Make sure to work on a copy of the spectra matrix, as that is much faster than row-wise extracting and changing the spectra by `[[` and `[[<-`.

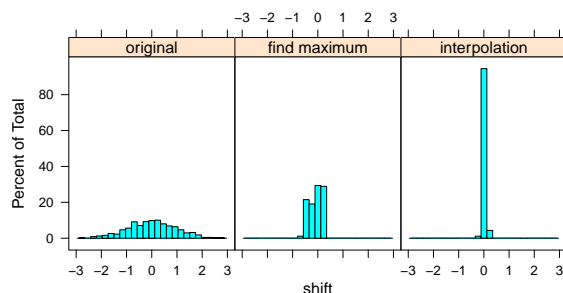


Figure 4: The shifts used to disturb the chondrocyte data (original), and the remaining shift after correction with the two methods discussed here.

```
> chondro <- chondro - spc.fit.poly.below (chondro [,min+3i ~ max - 3i], chondro)

Fitting with npts.min = 15

> chondro <- sweep (chondro, 1, rowMeans (chondro [[]], na.rm = TRUE), "/")

> targetfn <- function (shift, wl, spc, targetspc){
+   error <- spline (wl + shift, spc, xout = wl)$y - targetspc
+   sum (error^2)
+ }
> shift2 <- numeric (nrow (chondro))
> tmp <- chondro [[]]
> target <- colMeans (chondro [[]])
> for (i in 1 : nrow (chondro))
+   shift2 [i] <- unlist (optimize (targetfn, interval = c (-5, 5), wl = chondro@wavelength,
+                                 spc = tmp[i,], targetspc = target)$minimum)
```

Figure 4 shows that the second correction method works better for the chondrocyte data. This was expected, as the spectra are hardly or not oversampled, but are very similar to each other.

### 12.3. Smoothing Interpolation

Spectra acquired by grating instruments are frequently interpolated onto a new wavelength axis, e.g. because the unequal data point spacing should be removed. Also, the spectra can be smoothed: reducing the spectral resolution allows to increase the signal to noise ratio. For chemometric data analysis reducing the number of data points per spectrum may be crucial as it reduces the dimensionality of the data.

*hyperSpec* provides two functions to do so: `spc.bin` and `spc.loess`.

`spc.bin` bins the spectral axis by averaging every *by* data points.

```
> plot (paracetamol, wl.range = c (300 ~ 1800, 2800 ~ max), xoffset = 850)
> p <- spc.loess (paracetamol, c(seq (300, 1800, 2), seq (2850, 3150, 2)))
> plot (p, wl.range = c (300 ~ 1800, 2800 ~ max), xoffset = 850, col = "red", add = TRUE)
> b <- spc.bin (paracetamol, 4)
> plot (b, wl.range = c (300 ~ 1800, 2800 ~ max), xoffset = 850,
+       lines.args = list (pch = 20, cex = .3, type = "p"), col = "blue", add = TRUE)
```

`spc.loess` applies R's `loess` function for spectral interpolation. Figure 5 shows the result of interpolating from 300 to 1800 and 2850 to 3150  $\text{cm}^{-1}$  with 2  $\text{cm}^{-1}$  data point distance. This corresponds to a spectral resolution of about 4  $\text{cm}^{-1}$ , and the decrease in spectral resolution can be seen at the

`spc.bin`  
`spc.loess`

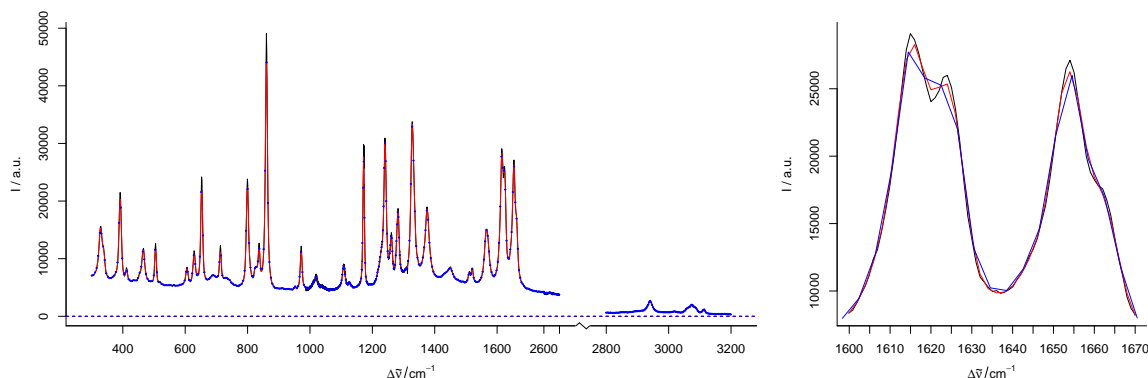


Figure 5: Smoothing interpolation by `spc.loess` with new data point spacing of  $2 \text{ cm}^{-1}$  (red) and `spc.bin` (blue). The magnification on the right shows how interpolation may cause a loss in signal height.

sharp bands where the maxima are not reached (due to the fact that the interpolation wavelength axis does not necessarily hit the maxima. The original spectrum had 4064 data points with unequal data point spacing (between 0 and  $1.4 \text{ cm}^{-1}$ ). The interpolated spectrum has 902 data points.

## 12.4. Background Correction

sweep

To subtract a background spectrum of each of the spectra in an object, use `sweep (spectra, 2, background.spectrum, "-")`.

## 12.5. Offset Correction

apply sweep

Calculate the offsets and sweep them off the spectra:

```
> offsets <- apply (chondro, 1, min)
> chondro.offset.corrected <- sweep (chondro, 1, offsets, "-")
```

If the offset is calculated by a function, as here with the `min`, *hyperSpec*'s `sweep` method offers a shortcut: `sweep`'s `STATS` argument may be the function instead of a numeric vector:

```
> chondro.offset.corrected <- sweep (chondro, 1, min, "-")
```

## 12.6. Baseline Correction

*hyperSpec* comes with two functions to fit polynomial baselines.

spc.fit.poly  
spc.fit.poly.below

`spc.fit.poly` fits a polynomial baseline of the given order. A least-squares fit is done so that the function may be used on rather noisy spectra. However, the user must supply an object that is cut appropriately. Particularly, the supplied wavelength ranges are not weighted.

`spc.fit.poly.below` tries to find appropriate support points for the baseline iteratively.

Both functions return a *hyperSpec* object containing the fitted baselines. They need to be subtracted afterwards:

```
> bl <- spc.fit.poly.below (chondro)
```

Fitting with `npts.min = 15`

```
> chondro <- chondro - bl
```

For details, see `vignette (baselinebelow)`.

## 12.7. Intensity Calibration

### 12.7.1. Correcting by a constant, e. g. Readout Bias

CCD cameras often operate with a bias, causing a constant value for each pixel. Such a constant can be immediately subtracted:

```
spectra - constant
```

### 12.7.2. Correcting Wavelength Dependence

sweep

For each of the wavelengths the same correction needs to be applied to all spectra.

1. There might be wavelength dependent offsets (background or dark spectra). They are subtracted:

```
sweep (spectra, 2, offset.spectrum, "-")
```

2. A multiplicative dependency such as a CCD's photon efficiency:

```
sweep (spectra, 2, photon.efficiency, "/")
```

### 12.7.3. Spectra Dependent Correction

sweep

If the correction depends on the spectra (e. g. due to inhomogeneous illumination while collecting imaging data, differing optical path length, etc.), the *MARGIN* of the `sweep` function needs to be 1 or SPC:

1. Pixel dependent offsets are subtracted:

```
sweep (spectra, SPC, pixel.offsets, "-")
```

2. A multiplicative dependency:

```
sweep (spectra, SPC, illumination.factors, "*")
```

## 12.8. Normalization

apply sweep

Again, `sweep` is the function of choice. E. g. for area normalization, use:

```
> chondro <- sweep (chondro, 1, mean, "/")
```

(using the mean instead of the sum results in conveniently scaled spectra with intensities around 1.)

If the calculation of the normalization factors is more elaborate, use a two step procedure:

1. Calculate appropriate normalization factors

You may calculate the factors using only a certain wavelength range, thereby normalizing on a particular band or peak.

2. Again, sweep the factor off the spectra:

```
normalized <- sweep (spectra, 1, factors, "*")
```

```
> factors <- 1 / apply (chondro [, , 1600 ~ 1700], 1, mean)
```

```
> chondro <- sweep (chondro, 1, factors, "*")
```

For minimum-maximum-normalization, first do an offset- or baseline correction, then normalize using `max`.

## 12.9. Centering the Data

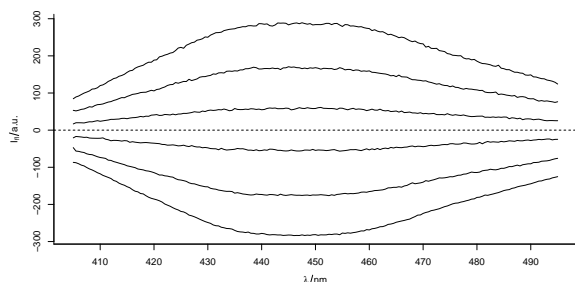
apply sweep

Centering means that the mean spectrum is subtracted from each of the spectra. Many data analysis techniques, like principal component analysis, partial least squares, etc., work much better on centered data.

However, from a spectroscopic point of view it depends on the particular data set whether centering does make sense or not.

To centre the `flu` data set, use:

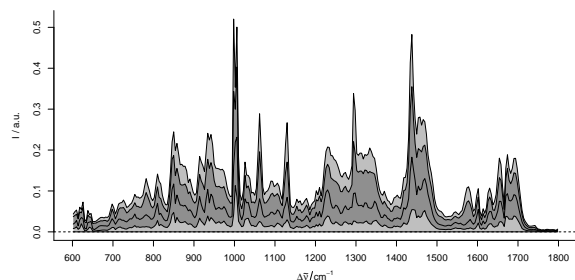
```
> flu.centered <- sweep (flu, 2, mean, "-")
> plot (flu.centered)
```



On the other hand, the `chondro` data set consists of Raman spectra, so the spectroscopic interpretation of centering is getting rid of the the average chemical composition of the sample. But: what is the meaning of the “average spectrum” of an inhomogeneous sample? In this case it may be better to subtract the minimum spectrum (which will hopefully have almost the same benefit on the data analysis) as it is the spectrum of that chemical composition that is underlying the whole sample.

One more point to consider is that the actual minimum spectrum will pick up (negative) noise. In order to avoid that, using e.g. the 5<sup>th</sup> percentile spectrum is more suitable:

```
> perc.5th <- apply (chondro, 2, quantile, 0.05)
> chondro <- sweep (chondro, 2, perc.5th, "-")
> plot (chondro, "spcprctl5")
```



## 12.10. Variance Scaling

apply sweep  
scale

Variance scaling is often used in multivariate analysis to adjust the influence and scaling of the variates (that are typically different physical values). However, spectra already do have the same scale of the same physical value. Thus one has to trade off the the expected numeric benefit with the fact that wavelengths with low signal will contain exploded noise after variance scaling.

Again, `sweep` may be used:

```
> scaled.chondro <- sweep (chondro, 2, var, "/")
```



Alternatively, R provides a function `scale` which works on matrices:

```
> scaled.chondro <- chondro
> scaled.chondro [[]] <- scale (scaled.chondro [[]])
```

### 12.11. Multiplicative Scatter Correction (MSC)

pls::msc

MSC can be done using `msc` from package *pls*[1]. It operates on the spectra matrix:

```
> require (pls)
> chondro.msc <- chondro
> chondro.msc [[]] <- msc (chondro [[]])
```

### 12.12. Spectral Arithmetic

+ - \* / ^ log  
log10

Basic mathematical functions are defined for *hyperSpec* objects. You may convert spectra:  
`absorbance.spectra = - log10 (transmission.spectra)`

In this case, do not forget to adapt the label:

labels

```
> labels (absorbance.spectra)$spc <- "A"
```

Be careful: R's `log` function calculates the natural logarithm if no base is given.

The basic arithmetic operators work element-wise in R. Thus they all need either a scalar, or a matrix (or *hyperSpec* object) of the correct size.

Matrix multiplication is done by `%*%`, again each of the operands may be a matrix or a *hyperSpec* object, and must have the correct dimensions. `%*%`

## 13. Data Analysis

### 13.1. Data Analysis Methods using a data.frame

#### e.g. Principal Component Analysis with `prcomp`

\$.

The `$.` notation is handy, if a data analysis function expects a *data.frame*. The column names can then be used in the formula:

```
> pca <- prcomp (~ spc, data = chondro$., center = FALSE)
```

Results of such a decomposition can be put again into *hyperSpec* objects. This allows to plot e.g. the loading like spectra, or score maps, see figure 6. `decomposition`

```
> scores <- decomposition (chondro, pca$x, label.wavelength = "PC",
+                           label.spc = "score / a.u.")
> scores
```

```
hyperSpec object
  875 spectra
   4 data columns
  300 data points / spectrum
wavelength: PC [integer] 1 2 ... 300
data: (875 rows x 4 columns)
  1. y: y/(mu * m) [numeric] -4.77 -4.77 ... 19.23
  2. x: x/(mu * m) [numeric] -11.55 -10.55 ... 22.45
  3. clusters: clusters [factor] matrix matrix ... lacuna + NA
  4. spc: score / a.u. [AsIs matrix x 300] -1.1696 -1.3931 ... -8.8349e-17
```

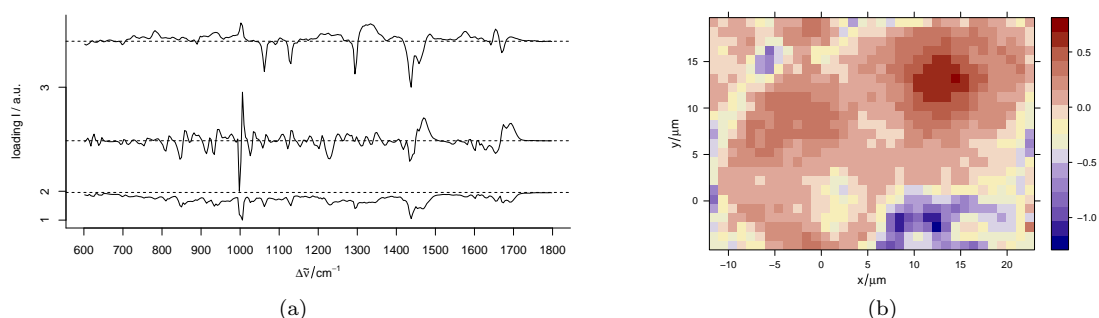


Figure 6: (a) The first three loadings: `plot (loadings [1 : 3], stacked = TRUE)`. (b) The third score map: `plotmap (scores [, , 3])`.

The loadings can be similarly re-imported:

```
> loadings <- decomposition (chondro, t(pca$rotation), scores = FALSE,
+                             label.spc = "loading I / a.u.")
> loadings

hyperSpec object
  300 spectra
  1 data columns
  300 data points / spectrum
wavelength: Delta * tilde(nu)/cm^-1 [numeric] 602 606 ... 1798
data: (300 rows x 1 columns)
  1. spc: loading I / a.u. [AsIs matrix x 300] -0.017699 -0.008177 ... 0.052487
```

There is, however, one important difference. The loadings are thought of as values computed from all spectra together. Thus no meaningful extra data can be assigned for the loadings object (at least not if the column consists of different values). Therefore, the loadings object lost all extra data (see above).

`retain.columns` triggers whether columns that contain different values should be dropped. If it is set to `TRUE`, the columns are retained, but contain NAs:

```
> loadings <- decomposition (chondro, t(pca$rotation), scores = FALSE,
+                             retain.columns = TRUE, label.spc = "loading I / a.u.")
> loadings[1]$..

  y x clusters
1 NA NA      <NA>
```

If an extra data column does contain only one unique value, it is retained anyways:

```
> chondro$measurement <- 1
> loadings <- decomposition (chondro, t(pca$rotation), scores = FALSE,
+                             label.spc = "loading I / a.u.")
> loadings[1]$..

  measurement
1            1
```

### 13.1.1. PCA as Noise Filter

Principal component analysis is sometimes used as a noise filtering technique. The idea is that the relevant differences are captured in the first components while the higher components contain noise only. Thus the spectra are reconstructed using only the first  $p$  components.

This reconstruction is in fact a matrix multiplication:

$$spectra^{(nrow \times nwl)} = scores^{(nrow \times p)} loadings^{(p \times nwl)}$$

Note that this corresponds to a model based on the Beer-Lambert law:

$$A_n(\lambda) = c_{n,i} \epsilon(i, \lambda) + error$$

The matrix formulation puts the  $n$  spectra into the rows of  $A$  and  $c$ , while the  $i$  pure components appear in the columns of  $c$  and rows of the absorbance coefficients  $\epsilon$ .

For an ideal data set (constituents varying independently, sufficient signal to noise ratio) one would expect the principal component analysis to extract something like the concentrations and pure component spectra.

If we decide that only the first 10 components actually carry spectroscopic information, we can reconstruct spectra with better signal to noise ratio:

```
> smoothed <- scores[, , 1:10] %*% loadings[1:10]
```

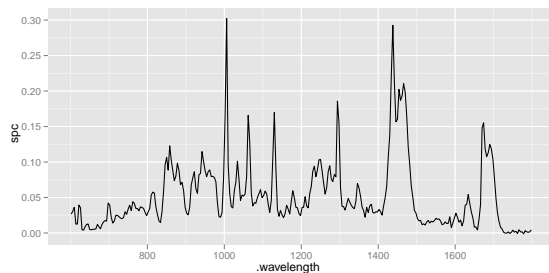
Keep in mind, though, that we cannot be sure how much *useful* information was discarded with the higher components. This kind of noise reduction may influence further modeling of the data. Mathematically speaking, the rank of the new  $875 \times 300$  spectra matrix is only 10.

## 13.2. Data Analysis using long-format data.frame

### e. g. plotting with ggplot2

Some functions need the data being an *unstacked* or *long-format data.frame*. `as.long.df` is the appropriate conversion function.

```
> require (ggplot2)
> ggplot (as.long.df (chondro [1]), aes (x = .wavelength, y = spc)) + geom_line ()
```



## 13.3. Data Analysis Methods using a matrix

### e. g. Hierarchical Cluster Analysis

```
> dist <- pearson.dist (chondro [[]])
> dendrogram <- hclust (dist, method = "ward")
> plot (dendrogram)
```

In order to plot a cluster map, the cluster membership needs to be calculated from the dendrogram. First, cut the dendrogram so that three clusters result:

```
> chondro$clusters <- as.factor (cutree (dendrogram, k = 3))
```

As the cluster membership was stored as factor, the levels can be meaningful names, which are displayed in the color legend.

```
> levels (chondro$clusters) <- c ("matrix", "lacuna", "cell")
```

Then the result may be plotted (figure 7b):

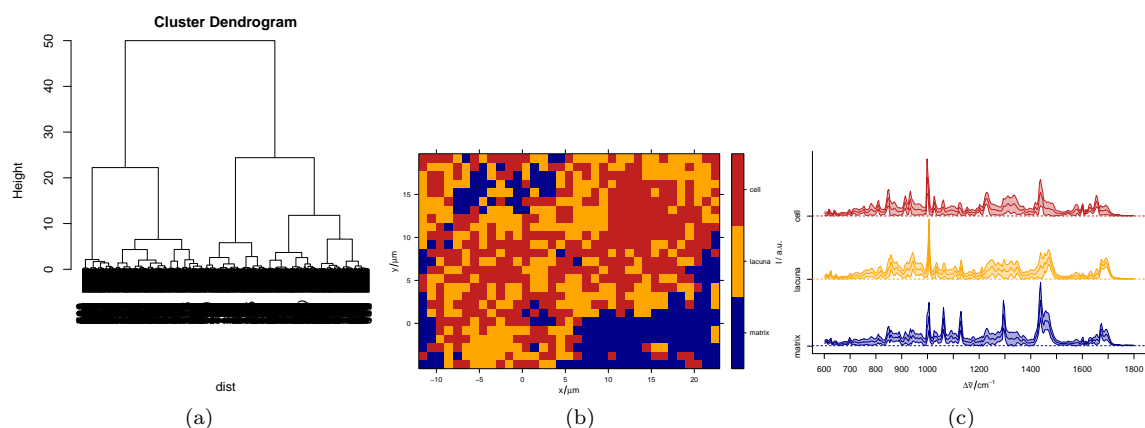


Figure 7: The results of the cluster analysis: (a) the dendrogram (b) the map of the 3 clusters (c) the mean spectra.

### 13.4. Calculating group-wise Sum Characteristics

#### e.g. Cluster Mean Spectra

`aggregate` applies the function given in *FUN* to each of the groups of spectra specified in *by*.

`aggregate`

So we may plot the cluster mean spectra:

```
> means <- aggregate (chondro, by = chondro$clusters, mean_pm_sd)
> plot (means, col = cluster.cols, stacked = ".aggregate", fill = ".aggregate")
```

### 13.5. Splitting an Object, and Binding a List of hyperSpec Objects

`split`

A *hyperSpec* object may also be split into a list of *hyperSpec* objects:

```
> clusters <- split (chondro, chondro$clusters)
> clusters

$matrix
hyperSpec object
  192 spectra
   5 data columns
 300 data points / spectrum
wavelength: Delta * tilde(nu)/cm^-1 [numeric] 602 606 ... 1798
data: (192 rows x 5 columns)
  1. y: y/(mu * m) [numeric] -4.77 -4.77 ... 19.23
  2. x: x/(mu * m) [numeric] -11.55 -10.55 ... 4.45
  3. clusters: clusters [factor] matrix matrix ... matrix
  4. spc: I / a.u. [matrix300] 0.026646 0.023320 ... 0.0032353
  5. measurement: measurement [numeric] 1 1 ... 1

$lacuna
hyperSpec object
  310 spectra
   5 data columns
 300 data points / spectrum
wavelength: Delta * tilde(nu)/cm^-1 [numeric] 602 606 ... 1798
data: (310 rows x 5 columns)
  1. y: y/(mu * m) [numeric] -4.77 -4.77 ... 19.23
  2. x: x/(mu * m) [numeric] -8.55 -5.55 ... 20.45
  3. clusters: clusters [factor] lacuna lacuna ... lacuna
```

```

4. spc: I / a.u. [matrix300] 0.048268 0.018306 ... 0.0009074
5. measurement: measurement [numeric] 1 1 ... 1

$cell
hyperSpec object
  373 spectra
  5 data columns
  300 data points / spectrum
wavelength: Delta * tilde(nu)/cm^-1 [numeric] 602 606 ... 1798
data: (373 rows x 5 columns)
  1. y: y/(mu * m) [numeric] -4.77 -4.77 ... 19.23
  2. x: x/(mu * m) [numeric] -7.55 -6.55 ... 22.45
  3. clusters: clusters [factor] cell cell ... cell
  4. spc: I / a.u. [matrix300] 0.028936 0.053938 ... 0.00033915
  5. measurement: measurement [numeric] 1 1 ... 1

```

Splitting can be reversed by `rbind` (see section 9.1, page 8). Another, similar way to combine a number of *hyperSpec* objects with different wavelength axes or extra data columns is `collapse` (see section 9.2, page 9).

## 14. Speed Considerations

While most of *hyperSpec*'s functions work at a decent speed for interactive sessions (of course depending on the size of the object), iterated (repeated) calculations as for bootstrapping or iterated cross validation may ask for special speed considerations.

As an example, let's again consider the code for shifting the spectra:

```

> tmp <- chondro [1 : 50]
> shifts <- rnorm (nrow (tmp))
> system.time ({
+   for (i in seq_len (nrow (tmp)))
+     tmp [[i]] <- interpolate (tmp [[i]], shifts [i], wl = wl (tmp))
+ })

      user system elapsed
    0.410   0.000   0.419

```

A first possibility to speed up is switching of the automatic logging of how the objects are transformed. Logging involves appending rows to the *data.frame* in slot `@log`. While the absolute amount of time needed to add a logbook entry is small, it may be executed very often (e.g. during each call of []).

```

> hy.setOptions (log = FALSE)
> tmp <- chondro [1 : 50]
> system.time ({
+   for (i in seq_len (nrow (tmp)))
+     tmp [[i]] <- interpolate (tmp [[i]], shifts [i], wl = wl (tmp))
+ })

      user system elapsed
    0.220   0.000   0.215

```

```

> hy.setOptions (log = TRUE)

```

Calculations that involve a lot of subsetting (i.e. extracting or changing the spectra matrix or extra data) can be sped up considerably if the required parts of the *hyperSpec* object are extracted beforehand. A related tip is that many model fitting functions in R are much faster if the formula interface is avoided and the appropriate *data.frames* or matrices are handed over directly.

```

> tmp <- chondro [1 : 50]
> system.time ({
+   tmp.matrix <- tmp [[]]
+   wl <- wl (tmp)
+   for (i in seq_len (nrow (tmp)))
+     tmp.matrix [i, ] <- interpolate (tmp.matrix [i, ], shifts [i], wl = wl)
+   tmp [[]] <- tmp.matrix
+ })

      user  system elapsed
0.020   0.000   0.024

```

## References

- [1] Ron Wehrens and Bjørn-Helge Mevik. *pls: Partial Least Squares Regression (PLSR) and Principal Component Regression (PCR)*, 2007. URL <http://mevik.net/work/software/pls.html>. R package version 2.1-0.

## A. Overview of the functions provided by hyperSpec

Function	Explanation
<i>Access parts of the object</i>	
[	Select / extract / delete spectra, wavelength ranges or extra data
[<-	Set parts of spectra or extra data
[[]	Select / extract / delete spectra, wavelength ranges or extra data, get the result as matrix or data.frame
[[]<-	Set parts of spectra matrix
\$	extract a data column (including \$spc)
\$<-	replace a data column (including \$spc)
i2wl	convert spectra matrix column indices to wavelengths
isample	get a random sample of the spectra as index vector
labels	get column labels
labels<-	set column labels
logbook	logging the data treatment
logentry	make a logbook entry
rownames<-	
sample	generate random sample of the spectra
seq.hyperSpec	sequence along the spectra, either as <i>hyperSpec</i> object or index vector
wl	extract the wavelengths
wl<-	replace the wavelengths
wl2i	convert wavelengths to spectra matrix column indices
<i>Basic information</i>	

Function	Explanation
<code>colnames</code>	
<code>colnames&lt;-</code>	
<code>dim</code>	
<code>dimnames</code>	
<code>length</code>	
<code>ncol</code>	number of data columns (extra data plus spectra matrix)
<code>nrow</code>	number of spectra
<code>nwl</code>	number of data points per spectrum
<code>print</code>	summary information
<code>rownames</code>	
<code>show</code>	
<code>summary</code>	summary information including the log
<code>chk.hy</code>	checks whether the object is a <code>hyperSpec</code> object
<i>Combine/split</i>	
<code>bind</code>	commom interface for <code>rbind</code> and <code>cbind</code>
<code>cbind2</code>	bind two <i>hyperSpec</i> objects by column
<code>cbind.hyperSpec</code>	
<code>collapse</code>	combine objects by adding columns if necessary. See <code>plyr::rbind.fill</code> .
<code>rbind2</code>	bind two <i>hyperSpec</i> objects by row, i. e. add wavelength ranges or extra data
<code>rbind.hyperSpec</code>	bind objects by row, i. e. add wavelength ranges or extra data
<code>split</code>	
<code>merge</code>	combines spectral ranges. works if spectra are in only one of the data sets
<i>Comparison</i>	
<code>all.equal</code>	
<code>Compare</code>	<code>&gt; &lt; == &gt;= &lt;=</code> return a logical matrix
<code>is.na</code>	
<i>Create and initialize an object</i>	
<code>initialize</code>	
<i>File import/export</i>	
<code>read.ENVI</code>	import ENVI file
<code>read.ENVI.Nicolet</code>	import ENVI files written by Nicolet spectrometers
<code>read.spc</code>	import .spc file
<code>read.spc.KaiserMap</code>	import a Raman map saved by Kaiser Optical Systems' Hologram software as multiple .spc files

Function	Explanation
<code>read.txt.long</code>	import long-type ASCII file
<code>read.txt.wide</code>	import wide-type ASCII file
<code>R.matlab::readMat</code>	import matlab file
<code>R.matlab::writeMat</code>	export as matlab file
<code>scan.txt.Renishaw</code>	import ASCII files produced by Renishaw (InVia) spectrometers
<code>write.txt.long</code>	export as long-type ASCII file
<code>write.txt.wide</code>	export as wide-type ASCII file
<code>scan.zip.Renishaw</code>	directly read zip packed ASCII files produced by Renishaw spectrometers
<i>Maths</i>	
<code>%*%</code>	matrix multiplication
<code>Arith</code>	
<code>log</code>	
<code>Math</code>	mathematical functions. See ( <code>help ("Math extquotedbl )</code> )
<code>Math2</code>	rounding
<code>Summary</code>	summary measures such as <code>min</code> , <code>max</code> , etc.
<i>Plotting</i>	
<code>levelplot</code>	
<code>map.identify</code>	identify spectra in map plot
<code>matlab.dark.palette</code>	darker version of <code>matlab.palette</code>
<code>matlab.palette</code>	palette resembling Matlab's jet colors
<code>plot</code>	main switchyard for plotting
<code>plotc</code>	intensity over one other dimension: calibration plots, time series, depth series, etc.
<code>plotmap</code>	false-colour intensity over two other dimensions: spectral images, maps, etc. (rectangular tessellation)
<code>plotspc</code>	spectra plots: intensity over wavelength
<code>plotvoronoi</code>	false-colour intensity over two other dimensions: spectral images, maps, etc. (Voronoi tessellation)
<code>spc.identify</code>	identify spectra and wavelengths in spectra plot
<code>stacked.offsets</code>	calculate intensity axis offsets for stacked spectral plots
<code>trellis.factor.key</code>	modify list of <code>levelplot</code> arguments according to factor levels
<i>Spectra-specific transformations</i>	
<code>orderwl</code>	sort columns of spectra matrix according to the wavelengths
<code>spc.bin</code>	spectral binning
<code>spc.fit.poly</code>	least squares fit of a polynomial
<code>spc.fit.poly.below</code>	least squares fit of a polynomial with automatic support point determination



Function	Explanation
<code>spc.loess</code>	loess smoothing interpolation
<i>Type conversion</i>	
<code>as.character</code>	
<code>as.data.frame</code>	
<code>as.long.df</code>	convert to a long-format data.frame.
<code>as.matrix</code>	
<code>as.wide.df</code>	convert to a wide-format data.frame with each wavelength one column
<code>decomposition</code>	re-import results of spectral matrix decomposition (or the like) into <i>hyperSpec</i> object
<i>Utility functions</i>	
<code>array2df</code>	convert array into a matrix or data.frame
<code>array2vec</code>	convert array indices ( $n$ element vector) into vector indices
<code>mean_pm_sd</code>	mean $\pm$ one standard deviation of a vector
<code>mean_sd</code>	mean and standard deviation of a vector
<code>pearson.dist</code>	distance measure based on Pearson's $R^2$
<code>rbind.fill</code>	transitional patch of <code>plyr::rbind.fill</code> working with matrices
<code>rbind.fill.matrix</code>	transitional until <code>plyr::rbind.fill.matrix</code> is out
<code>vec2array</code>	convert vector (one element) index into an array into an $n$ element array index
<code>wc</code>	word count using <code>wc</code> if available on the system
<i>Vectorization</i>	
<code>aggregate</code>	
<code>apply</code>	
<code>sweep</code>	

## Session Info

R version 2.12.1 (2010-12-16)

Platform: x86\_64-pc-linux-gnu (64-bit)

locale:

```
[1] LC_CTYPE=en_US.utf8      LC_NUMERIC=C              LC_TIME=en_US.utf8
[4] LC_COLLATE=en_US.utf8    LC_MONETARY=C             LC_MESSAGES=en_US.utf8
[7] LC_PAPER=en_US.utf8      LC_NAME=C                  LC_ADDRESS=C
[10] LC_TELEPHONE=C           LC_MEASUREMENT=en_US.utf8 LC_IDENTIFICATION=C
```

attached base packages:

```
[1] grid      stats      graphics  grDevices  utils      datasets  methods   base
```

other attached packages:

```
[1] ggplot2_0.8.9      proto_0.3-8      reshape_0.8.4      plyr_1.4
[5] plotrix_3.0-8      hyperSpec_0.96-20110218 lattice_0.19-17
```