

# IPPD package vignette

Martin Slawski <sup>1</sup>	<a href="mailto:ms@cs.uni-saarland.de">ms@cs.uni-saarland.de</a>
Rene Hussong <sup>2</sup>	<a href="mailto:rene.hussong@uni.lu">rene.hussong@uni.lu</a>
Andreas Hildebrandt <sup>3</sup>	<a href="mailto:andreas.hildebrandt@uni-mainz.de">andreas.hildebrandt@uni-mainz.de</a>
Matthias Hein <sup>1</sup>	<a href="mailto:hein@cs.uni-saarland.de">hein@cs.uni-saarland.de</a>

<sup>1</sup> Saarland University, Department of Computer Science, Machine Learning Group, Saarbrücken, Germany

<sup>2</sup> Luxembourg Centre for Systems Biomedicine, University of Luxembourg, Luxembourg

<sup>3</sup> Johannes Gutenberg-University Mainz, Institute for Computer Science, Mainz, Germany

,

## Abstract

This is the vignette of the Bioconductor add-on package IPPD which implements automatic isotopic pattern extraction from a raw protein mass spectrum. Basically, the user only has to provide mass/charge channels and corresponding intensities, which are automatically decomposed into a list of monoisotopic peaks. IPPD can handle several charge states as well as overlaps of peak patterns.

## 1 Aims and scope of IPPD

A crucial challenge in the analysis of protein mass spectrometry data is to automatically process the raw spectrum to a list of peptide masses. IPPD is tailored to spectra where peptides emerge in the form of isotope patterns, i.e. one observes several peaks for each peptide mass at a given charge state due to the natural abundance of heavy isotopes. Datasets with a size of up to 100,000 mass/charge channels and the presence of isotope patterns at multiple charge states frequently exhibiting overlap make the manual annotation of a raw spectrum a tedious task. IPPD provides functionality to perform this task in a fully automatic, transparent and user-customizable way. Basically, one feeds the raw spectrum into one single function to obtain a list of monoisotopic peaks described by a mass/charge channel, a charge and an intensity. What makes our approach particularly user-friendly is its dependence on only a small set of easily interpretable parameters. We also offer a method to display the decomposition of the spectrum graphically, thereby facilitating a manual validation of the output.

## 2 Methodology

### 2.1 Template model

In the context of this package, a protein mass spectrum is understood as a sequence of pairs  $\{x_i, y_i\}_{i=1}^n$ , where  $x_i = m_i/z_i$  is a mass ( $m_i$ ) per charge ( $z_i$ ) value (measured in Thomson) and  $y_i$  is the intensity, i.e. the abundance of a particular mass (modulo charge state),



Figure 1: Illustration of the template construction as described in the text. The left panel depicts different templates of different charge states (1 to 4). The right panel zooms at the charge two template  $\varphi_2$ .

observed at  $x_i$ ,  $i = 1, \dots, n$ , which are assumed to be in an increasing order. The  $y_i$  are modeled as a linear combination of template functions representing prior knowledge about peak shapes and the composition of isotopic patterns. If our model were exact, we could write

$$\mathbf{y} = \Phi \boldsymbol{\beta}^*, \quad \mathbf{y} = (y_1, \dots, y_n)^\top, \quad (1)$$

where  $\Phi$  is a matrix template functions and  $\boldsymbol{\beta}^*$  a vector of weights for each template. Only a small fraction of all templates are needed to fit the signal, i.e.  $\boldsymbol{\beta}^*$  is highly sparse. Since  $\mathbf{y} \geq 0$ , where ' $\geq$ ' is understood componentwise, all template functions are nonnegative and accordingly  $\boldsymbol{\beta}^* \geq 0$ . Model (1) can equivalently be written as

$$\mathbf{y} = \begin{bmatrix} \Phi_1 & \dots & \Phi_C \end{bmatrix} \begin{bmatrix} \beta_1^* \\ \vdots \\ \beta_C^* \end{bmatrix} = \sum_{c=1}^C \Phi_c \beta_c^*, \quad (2)$$

where  $\Phi_c, \beta_c^*$  denote the matrix of template functions and weight vector to fit isotopic patterns of a particular charge state  $c$ ,  $c = 1, \dots, C$ . Each submatrix  $\Phi_c$  can in turn be divided into columns  $\varphi_{c,1}, \dots, \varphi_{c,p_c}$ , where the entries of each column vector store the evaluations of a template  $\varphi_{c,j}$ ,  $j = 1, \dots, p_c$ , at the  $x_i$ ,  $i = 1, \dots, n$ . Each template  $\varphi_{c,j}$  depends on parameter  $m_{c,j}$  describing the  $m/z$  position at which  $\varphi_{c,j}$  is placed. A template  $\varphi_{c,j}$  is used to fit an isotopic pattern of peaks composed of several single peaks, which is modeled as

$$\varphi_{c,j} = \sum_{k \in Z_{c,j}} a_{c,j,k} \psi_{c,j,k, \boldsymbol{\theta}_{c,j}}, \quad Z_{c,j} \subset \mathbb{Z} \quad (3)$$

where the  $\psi_{c,j,k}$  are functions representing a peak of a single isotope within an isotopic pattern. They depend on  $m_{c,j}$  and a parameter vector  $\boldsymbol{\theta}_{c,j}$ . The nonnegative weights  $a_{c,j,k}$  reflect the relative abundance of the isotope indexed by  $k$ . The  $a_{c,j,k}$  are computed

according to the averagine model (Senko et al. [1995]) and hence are fixed in advance. Each  $\psi_{c,j,k}$  is linked to a location  $m_{c,j,k}$  at which it attains its maximum. The  $m_{c,j,k}$  are calculated from  $m_{c,j}$  as  $m_{c,j,k} = m_{c,j} + \kappa \frac{k}{c}$ , where  $\kappa$  equals 1 Dalton ( $\approx 1.003$ ). The rationale behind Eq. (3) and the definitions that follow is the fact that the location of the most intense isotope is taken as characteristic location of the template, i.e. we set  $m_{c,j,0} = m_{c,j}$  so that the remaining  $m_{c,j,k}$ ,  $k \neq 0$ , are computed by shifting  $m_{c,j}$  in both directions on the  $m/z$  axis. By 'most intense isotope', we mean that  $a_{c,j,0} = \max_k a_{c,j,k} = 1$ . The set  $Z_{c,j}$  is a subset of the integers which depends on the averagine model and a pre-specified tolerance, i.e. we truncate summation in Eq. (3) if the weights drop below that tolerance. Figure 1 illustrates the construction scheme and visualizes our notation.

## 2.2 Peak shape

In an idealized setting, the  $\psi_{c,j,k}$  are delta functions at specific locations. In practice, however, the shape of a peak equals that of a bump which may exhibit some skewness. In the case of no to moderate skewness, we model peaks by Gaussian functions:

$$\psi_{c,j,k}(x) = \exp\left(-\frac{(x - m_{c,j,k})^2}{\sigma_{c,j}}\right). \quad (4)$$

The parameter to be determined is  $\theta_{c,j} = \sigma_{c,j} > 0$ . In the case of considerable skewness, peaks are modeled by exponentially modified Gaussian (EMG) functions, see for instance Grushka [1972], Marco and Bombi [2001], and Schulz-Trieglaff et al. [2007] in the context of protein mass spectrometry:

$$\begin{aligned} \psi_{c,j,k}(x) &= \frac{1}{\alpha_{c,j}} \exp\left(\frac{\sigma_{c,j}^2}{2\alpha_{c,j}^2} + \frac{\mu_{c,j} - (x - m_{c,j,k})}{\alpha_{c,j}}\right) \left(1 - F\left(\frac{\sigma_{c,j}}{\alpha_{c,j}} + \frac{\mu_{c,j} - (x - m_{c,j,k})}{\sigma_{c,j}}\right)\right), \\ F(t) &= \int_{-\infty}^t \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{u^2}{2}\right) du. \end{aligned} \quad (5)$$

The EMG function involves a vector of three parameters  $\theta_{c,j} = (\alpha_{c,j}, \sigma_{c,j}, \mu_{c,j})^\top \in \mathbb{R}^+ \times \mathbb{R}^+ \times \mathbb{R}$ . The parameter  $\alpha_{c,j}$  controls the additional length of the right tail as compared to a Gaussian. For  $\alpha_{c,j} \downarrow 0$ , the EMG function becomes a Gaussian. For our fitting approach as outlined in Section 2.3, it is crucial to estimate the  $\theta_{c,j}$ , which are usually unknown, from the data as good as possible. To this end, we model each component  $\theta_l$  of  $\theta$  as a linear combination of known functions  $g_{l,m}$  of  $x = m/z$  and an error component  $\varepsilon_l$ , i.e.

$$\theta_l(x) = \sum_{m=1}^{M_l} \nu_{l,m} g_{l,m}(x) + \varepsilon_l(x). \quad (6)$$

In the case of no prior knowledge about the  $g_{l,m}$ , we model  $\theta_l$  as a constant independent of  $x$ . In most cases, it is sensible to assume a linear trend, i.e.  $\theta_l(x) = \nu_{l,1} + \nu_{l,2}x$ . In order to fit a model of the form (6), we have to collect information from the data  $\{x_i, y_i\}_{i=1}^n$ . To be precise, we proceed according to the following steps.

1. We apply a simple peak detection algorithm to the spectrum to identify disjoint regions  $\mathcal{R}_r \subset \{1, \dots, n\}$ ,  $r = 1, \dots, R$ , of well-resolved peaks.

2. For each region  $r$ , we fit the chosen peak shape to the data  $\{x_i, y_i\}_{i \in \mathcal{R}_r}$  using nonlinear least squares:

$$\min_{\boldsymbol{\theta}} \sum_{i \in \mathcal{R}_r} (y_i - \psi_{\boldsymbol{\theta}}(x_i))^2, \quad (7)$$

yielding an estimate  $\hat{\boldsymbol{\theta}}_r(\hat{x}_r)$ , where  $\hat{x}_r$  denotes an estimation for the mode of the peak in region  $\mathcal{R}_r$ .

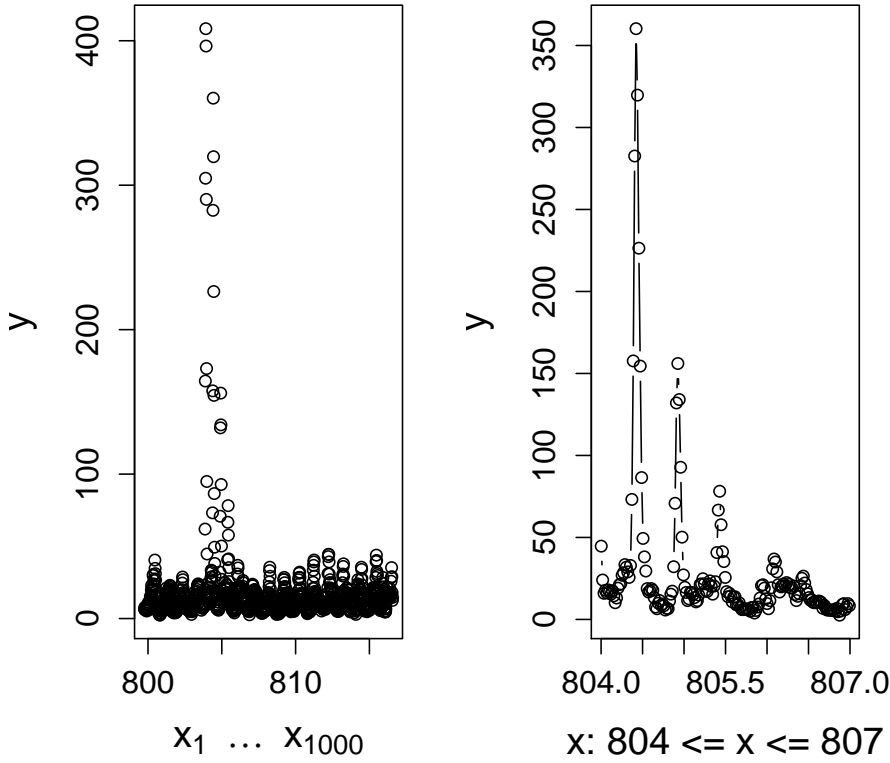
3. The sequence  $\{\hat{x}_r, \hat{\boldsymbol{\theta}}_r\}_{r=1}^R$  is then used as input for the estimation of the parameters  $\nu_{l,m}$  in model (6).

Step 2. is easily solved by the general purpose nonnegative least squares routine `nls` in `R:::stats` for a Gaussian peak shape. For the EMG, we have to perform a grid search over all three parameters to find a suitable starting value, which is then passed to the general purpose optimization routine `optim` in `R:::stats` with the option `method = "BFGS"` and a specification of a closed form expression of the gradient via the argument `gr`. For step 3., we use least absolute deviation regression because of the presence of outliers arising from less well-resolved, wiggly or overlapping peaks. The whole procedure is performed by the function `fitModelParameters` as demonstrated below. After loading the package, we access the real world dataset `myo500` and extract  $m/z$  channels (`x`) and the corresponding intensities (`y`). For computational convenience and since they contain very few relevant information, we discard all channels above 2500.

```
R> library(IPPD)
R> data(myo500)
R> x <- myo500[, "mz"]
R> y <- myo500[, "intensities"]
R> y <- y[x <= 2500]
R> x <- x[x <= 2500]
```

To have a look at the data, we plot the first 1000 (`x`,`y`) pairs:

```
R> layout(matrix(c(1,2), 1, 2))
R> plot(x[1:1000], y[1:1000], xlab = expression(x[1]~~ldots~~x[1000]),
       cex.lab = 1.5, cex.axis = 1.25, ylab = expression(y))
R> plot(x[x >= 804 & x <= 807], y[x >= 804 & x <= 807],
       xlab = "x: 804 <= x <= 807",
       cex.lab = 1.5, cex.axis = 1.25, ylab = expression(y), type = "b")
R> layout(matrix(1))
R>
```



In the plot, one identifies a prominent peak pattern beginning at about 804, which is zoomed at in the right panel.

We now apply `fitModelParameters` to fit model (6) for the width parameter  $\sigma$  of a Gaussian function (4). For simplicity, we take  $g_1(x) = 1$ ,  $g_2(x) = x$ . The model is specified by using an R formula interface.

```
R> fitGauss <- fitModelParameters(mz = x, intensities = y,
                                model = "Gaussian", fitting = "model", formula.sigma = formula(~mz),
                                control = list(window = 6, threshold = 200))
```

An analogous command for the EMG (5) with the model formulae  $\alpha(x) = \nu_{1,1} + \nu_{1,2}x$ ,  $\sigma(x) = \nu_{2,1} + \nu_{2,2}x$ ,  $\mu(x) = \nu_{3,1}$  is given by

```
R> fitEMG <- fitModelParameters(mz = x, intensities = y,
                                model = "EMG", fitting = "model",
                                formula.alpha = formula(~mz),
                                formula.sigma = formula(~mz),
                                formula.mu = formula(~1),
                                control = list(window = 6, threshold = 200))
```

Inspecting the results, we find that  $R = 55$  peak regions are used to fit an EMG parameter model. Moreover, it turns out that the EMG model is a more appropriate peak model for the data when visually comparing the list of mean residual sums of squares of the EMG fits and the Gauss fits extracted from `slot(fitEMG, "peakfitresults")` and

`slot(fitGauss, "peakfitresults")`, respectively. The figure shows an example where the EMG shape comes relatively close to the observed data. A long right tail indicates that a Gaussian would yield a rather poor fit here.

```
R> show(fitEMG)
```

```
Peak model 'EMG' fitted as fuction of m/z
number of peaks used: 55
```

```
R> mse.EMG <- data.frame(mse = slot(fitEMG,"peakfitresults")[,"rss"]
                        / slot(fitEMG,"peakfitresults")[,"datapoints"],
                        peakshape = rep("EMG", nrow( slot(fitEMG,"peakfitresults"))))
R> mse.Gauss <- data.frame(mse = slot(fitGauss,"peakfitresults")[,"rss"]
                        / slot(fitGauss,"peakfitresults")[,"datapoints"],
                        peakshape = rep("Gaussian", nrow( slot(fitGauss,"peakfitresults"))))
R> mses <- rbind(mse.EMG, mse.Gauss)
R> with(mses, boxplot(mse ~ peakshape, cex.axis = 1.5, cex.lab = 1.5, ylab = "MSE"))
R>
```

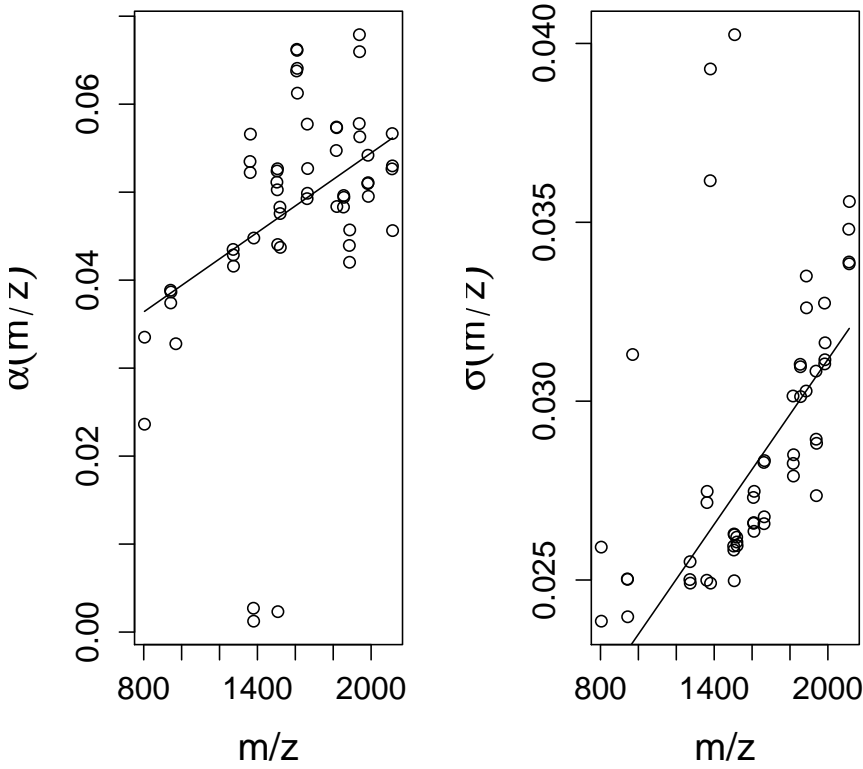


```
R> visualize(fitEMG, type = "peak", cex.lab = 1.5, cex.axis = 1.25)
```



To assess the fit of the two linear models for the EMG parameters  $\alpha$  and  $\sigma$ , we use again the function `visualize` as follows:

```
R> visualize(fitEMG, type = "model", modelfit = TRUE,
  parameters = c("sigma", "alpha"),
  cex.lab = 1.5, cex.axis = 1.25)
R>
```



While the fit for  $\sigma$  seems to be reasonable except for some extreme outliers, the fit for  $\alpha$  is not fully convincing. Nevertheless, in the absence of further knowledge, the fit produces good results in the template matching step detailed in the next section.

### 2.3 Template fitting

Once all necessary parameters have been determined, the positions at which the templates are placed have to be fixed. In general, one has to choose positions from the interval  $[x_1, x_n]$ . We instead restrict us to a suitable subset of the finite set  $\{x_i\}_{i=1}^n$ . The deviations from the true positions is then at least in the order of the sampling rate, but this can be improved by means of a postprocessing step described in 2.4. Using the whole set  $\{x_i\}_{i=1}^n$  may be computationally infeasible if  $n$  is large. Such an approach would be at least computationally wasteful, since 'genuine' peaks patterns occur very sparsely in the spectrum. Therefore, we apply a pre-selection step on the basis of what we term 'local noise level' (LNL). The LNL is defined as a quantile (typically the median) of the intensities  $y_i$  falling into a sliding window of fixed width around a specific position. Given the LNL, we place templates on an  $x_i$  (one for each charge state) if and only if the corresponding  $y_i$  exceeds the LNL at  $x_i$  by a factor `factor.place`, which typically equals three or four and has to be specified by the user. Given the positions of the templates, we compute the matrix  $\Phi$  according to Eqs. (1) and (3). It then remains to estimate the coefficient vector  $\beta^*$  on the basis of two structural assumptions, sparsity and nonnegativity of all quantities involved. Related approaches in the literature (Du and Angeletti [2006], Renard et al. [2008]) account for sparsity of  $\beta^*$  by using  $\ell_1$ -regularized regression (Tibshirani [1996]). We here argue empirically that  $\ell_1$  regularization is not the best to do, since it entails the



selection of a tuning parameter which is difficult to choose in our setting, and secondly the structural constraints concerning nonnegativity turn out to be so strong that sparsity is more conveniently achieved by fitting followed by hard thresholding. We first determine

$$\begin{aligned} \hat{\beta} \in \underset{\beta}{\operatorname{argmin}} \|\mathbf{y} - \Phi\beta\|_q^q, \quad q = 1 \text{ or } q = 2, \\ \text{subject to } \beta \geq 0. \end{aligned} \quad (8)$$

The optimization problem (8) is a quadratic ( $q = 2$ ) or linear ( $q = 1$ ) program and is solved using standard techniques (Boyd and Vandenberghe [2004]); we omit further details here. We remark that in the presence of high noise, it is helpful to subtract the LNL from  $\mathbf{y}$ . Concerning the choice of  $q$ , we point out that  $q = 1$  can cope better with deviations from model assumptions, i.e. deviations from the average model or from the peak model and thus may lead to a reduction of the number of false positives.

## 2.4 Postprocessing

Given an estimate  $\hat{\beta}$ , we define  $\mathcal{M}_c = \{m_{c,j} : \hat{\beta}_{c,j} > 0\} \subset \{x_i\}_{i=1}^n$ ,  $c = 1, \dots, C$ , as the set of all template locations where the corresponding coefficient exceeds 0, separately for each charge. Due to a limited sampling rate, different sources of noise and model misfit, the locations in the sets  $\{\mathcal{M}_c\}_{c=1}^C$  may still deviate considerably from the set of true peak pattern locations. Specifically, the sets  $\{\mathcal{M}_c\}_{c=1}^C$  tend to be too large, mainly caused by what we term 'peak splitting': for the reasons just mentioned, it frequently occurs that several templates are used to fit the same peak. This can at least partially be corrected by means of the following merging procedure.

1. Separately for each  $c$ , divide the sets  $\mathcal{M}_c$  into groups  $\mathcal{G}_{c,1}, \dots, \mathcal{G}_{c,G_c}$  of 'adjacent' positions. Positions are said to be adjacent if their distance on the  $m/z$  scale is below a certain tolerance as specified via a parts per million (ppm) value.
2. For each  $c = 1, \dots, C$  and each group  $g_c = 1, \dots, G_c$ , we solve the following optimization problem.

$$(\tilde{m}_{c,g}, \tilde{\beta}_{c,g}) = \min_{m_{c,g}, \beta_{c,g}} \left\| \sum_{m_{c,j} \in \mathcal{G}_{c,g}} \hat{\beta}_{c,j} \psi_{m_{c,j}} - \beta_{c,g} \psi_{m_{c,g}} \right\|_{L^2}^2 \quad (9)$$

In plain words, we take the fitted function resulting from the functions  $\{\psi_{m_{c,j}}\}$  representing the most intense peak of each peak pattern in the same group and then determine a function  $\psi_{\tilde{m}_{c,g}}$  placed at location  $\tilde{m}_{c,g}$  and weighted by  $\tilde{\beta}_{c,g}$  such that  $\tilde{\beta}_{c,g} \psi_{\tilde{m}_{c,g}}$  approximates the fit of multiple functions  $\{\psi_{m_{c,j}}\}$  best (in a least squares sense).

3. One ends up with sets  $\tilde{\mathcal{M}}_c = \{\tilde{m}_{c,g}\}_{g=1}^{G_c}$  and coefficients  $\{\tilde{\beta}_{c,g}\}_{g=1}^{G_c}$ ,  $c = 1, \dots, C$ .

The additional benefit of step 2. as compared to the selection of the function with the largest coefficient as proposed in Renard et al. [2008] is that, in the optimal case, we are able to determine the peak pattern location even more accurate as predetermined by a limited sampling rate. The integral in (9) can be solved analytically for a Gaussian function, and we resort to numeric approximations for the EMG function.

The sets  $\{\mathcal{M}_c\}$  tend to be too large in the sense that they still contain noise peak patterns.

Therefore, we apply hard thresholding to the  $\{\tilde{\beta}_{c,g}\}_{g=1}^{G_c}$ ,  $c = 1, \dots, C$ , discarding all positions where the corresponding coefficients is less than a significance level times the LNL, where the significance level has to be specified by the user.

### 3 Case study

We continue the data analysis starting in Section 2.2. The methodology of the Sections 2.3 and 2.4 is implemented in the function `getPeaklist`. For the computation of the template functions, we recycle the object `fitEMG` obtained in Section 2.2.

```
R> EMGlist <- getPeaklist(mz = x, intensities = y, model = "EMG",
  model.parameters = fitEMG,
  loss = "L2", trace = FALSE,
  control.localnoise = list(factor.place = 2),
  control.basis = list(charges = c(1, 2)),
  control.postprocessing = list(ppm = 200))
R> show(EMGlist)
```

```
An object of class 'peaklist'(with postprocessing)
Loss function used: L2
Peak model used: EMG
number of peaks: 1222
charge states used: 1,2
```

```
R>
```

The argument list can be summarized as follows: we compute EMG templates for charges 1 and 2; templates are placed on all  $m/z$ -positions in the spectrum where the intensity is at least two times the LNL; the fit is least squares (`loss = L2`); postprocessing is performed by merging peaks within a tolerance of 200 ppm. Subsequently, only the patterns with signal-to-noise ratio bigger than three are maintained. The result is of the following form.

```
R> threshold(EMGlist, threshold = 3, refit = TRUE, trace = FALSE)
```

	loc_init	loc_most_intense	charge	quant	amplitude	localnoise	ratio
[1,]	800.4642	800.4642	1	79.45389	50.281828	9.411770	5.342441
[2,]	808.2414	808.2414	1	80.37673	50.671012	10.196100	4.969646
[3,]	829.2292	829.2292	1	93.78790	58.518172	9.411770	6.217552
[4,]	842.4935	842.4935	1	65.96983	40.891818	8.366010	4.887852
[5,]	864.4065	864.4065	1	142.84891	87.579019	13.333300	6.568443
[6,]	877.0373	877.0373	1	115.99279	70.662695	8.888890	7.949552
[7,]	908.4247	908.4247	1	77.96504	46.750652	8.104580	5.768424
[8,]	908.9339	908.9339	1	57.32357	34.365178	9.411770	3.651298
[9,]	923.4380	923.4380	1	71.88843	42.791000	7.581700	5.643985
[10,]	924.4543	924.4543	1	72.76197	43.288705	8.366010	5.174355
[11,]	927.4746	927.4746	1	109.42451	65.004254	8.888890	7.312978
[12,]	927.9689	927.9689	1	60.26112	35.789855	9.934640	3.602532
[13,]	941.4672	941.4672	1	6892.94438	4066.651284	12.549000	324.061781

[14,]	941.6761	941.6761	1	136.02533	80.240039	13.594800	5.902260
[15,]	942.4249	942.4249	1	272.44888	160.650355	16.732000	9.601384
[16,]	949.4364	949.4364	1	71.82305	42.204271	7.320260	5.765406
[17,]	952.5318	952.5318	1	76.12937	44.665674	7.843140	5.694871
[18,]	954.4819	954.4819	1	106.30197	62.306879	8.104580	7.687860
[19,]	963.4553	963.4553	1	148.96394	86.918604	6.535950	13.298542
[20,]	967.4826	967.4826	1	145.36799	84.650178	8.104580	10.444733
[21,]	969.4694	969.4694	1	262.70495	152.821228	10.457500	14.613553
[22,]	985.4424	985.4424	1	125.58555	72.468389	7.320260	9.899702
[23,]	991.5025	991.5025	1	143.73336	82.684219	9.411770	8.785193
[24,]	992.0244	992.0244	1	126.08040	72.510535	10.196100	7.111595
[25,]	999.4665	999.4665	1	153.66239	88.037291	7.320260	12.026525
[26,]	1023.4509	1023.4509	1	72.98976	41.206802	6.535950	6.304639
[27,]	1057.4478	1057.4478	1	61.60008	34.042067	7.320260	4.650390
[28,]	1086.5573	1086.5573	1	163.03327	88.434884	6.013070	14.707110
[29,]	1125.5187	1125.5187	1	82.20092	43.611929	6.797390	6.415982
[30,]	1151.4789	1151.4789	1	97.87418	51.214203	5.751630	8.904294
[31,]	1168.6235	1168.6235	1	88.01377	45.630786	6.274510	7.272406
[32,]	1192.7011	1192.7011	1	73.30719	37.510877	6.535950	5.739162
[33,]	1271.6609	1271.6609	1	3373.99411	1646.830797	10.457500	157.478441
[34,]	1297.6800	1297.6800	1	143.24711	68.802828	7.581700	9.074855
[35,]	1360.7611	1360.7611	1	3817.19023	1720.274078	13.071900	131.600921
[36,]	1361.7379	1361.7379	1	1393.40685	627.304820	20.653600	30.372662
[37,]	1378.8379	1378.8379	1	2998.58306	1325.589299	15.686300	84.506180
[38,]	1394.8413	1394.8413	1	133.39966	57.963577	9.150330	6.334589
[39,]	1474.6387	1474.6387	1	154.37366	63.766703	9.934640	6.418622
[40,]	1484.6606	1484.6606	1	554.64563	227.682827	11.764700	19.353050
[41,]	1500.6588	1500.6588	1	187.02460	76.020197	14.902000	5.101342
[42,]	1501.6659	1501.6659	1	237.91115	96.654308	20.130700	4.801339
[43,]	1502.6668	1502.6668	1	6512.56220	2644.229261	27.973900	94.524870
[44,]	1506.9383	1506.9383	1	1362.38924	551.970349	18.039200	30.598383
[45,]	1518.6637	1518.6637	1	3192.43548	1285.329890	14.902000	86.252174
[46,]	1519.6113	1519.6113	1	167.74745	67.506458	15.424800	4.376488
[47,]	1524.6527	1524.6527	1	142.45911	57.178770	9.411770	6.075241
[48,]	1534.6603	1534.6603	1	271.31959	108.323839	12.287600	8.815704
[49,]	1546.6550	1546.6550	1	180.74626	71.701512	10.196100	7.032249
[50,]	1588.8538	1588.8538	1	206.84694	80.203379	12.287600	6.527180
[51,]	1589.8317	1589.8317	1	205.68650	79.707376	14.379100	5.543280
[52,]	1606.8601	1606.8601	1	47659.91282	18289.702125	27.451000	666.267244
[53,]	1622.8482	1622.8482	1	282.49197	107.343397	8.104580	13.244782
[54,]	1628.8462	1628.8462	1	197.47847	74.759773	7.843140	9.531868
[55,]	1632.8754	1632.8754	1	231.62212	87.464920	9.673200	9.041984
[56,]	1643.8448	1643.8448	1	295.36538	110.772870	10.980400	10.088236
[57,]	1650.8348	1650.8348	1	219.30843	81.886503	9.411770	8.700436
[58,]	1660.8520	1660.8520	1	346.30112	128.490258	15.424800	8.330109
[59,]	1661.8523	1661.8523	1	7081.99298	2625.671093	20.392200	128.758599
[60,]	1675.8106	1675.8106	1	140.11613	51.492824	7.058820	7.294820
[61,]	1683.8293	1683.8293	1	115.18493	42.112772	6.535950	6.443252
[62,]	1687.8674	1687.8674	1	138.47730	50.496544	7.581700	6.660319
[63,]	1712.6676	1712.6676	1	116.18272	41.773239	6.274510	6.657610

[64,]	1717.8060	1717.8060	1	100.80201	36.153297	6.274510	5.761931
[65,]	1718.8742	1718.8742	1	78.34668	28.084804	6.535950	4.296973
[66,]	1753.7102	1753.7102	1	89.35261	31.479802	5.490200	5.733817
[67,]	1777.8665	1777.8665	1	106.10370	36.931121	6.013070	6.141808
[68,]	1798.8819	1798.8819	1	175.36393	60.387426	7.581700	7.964893
[69,]	1815.9052	1815.9052	1	9222.23376	3149.090676	24.052300	130.926800
[70,]	1831.9004	1831.9004	1	327.93247	111.106409	5.751630	19.317378
[71,]	1837.8935	1837.8935	1	368.32791	124.420732	6.274510	19.829554
[72,]	1847.8954	1848.9034	1	373.94066	125.769264	6.274510	20.044476
[73,]	1852.9575	1853.9655	1	162.82771	54.759202	10.980400	4.986995
[74,]	1853.9657	1854.9737	1	2843.33423	956.199297	12.549000	76.197251
[75,]	1869.9595	1870.9675	1	428.17483	143.949541	6.013070	23.939442
[76,]	1885.0257	1886.0337	1	2133.58705	717.090660	9.150330	78.367738
[77,]	1885.9545	1886.9625	1	157.60489	52.969416	9.673200	5.475894
[78,]	1897.9419	1898.9499	1	73.78245	24.791914	4.444440	5.578186
[79,]	1901.0127	1902.0207	1	88.40198	29.699518	4.705880	6.311151
[80,]	1919.0078	1920.0158	1	146.81338	49.262011	6.013070	8.192489
[81,]	1937.0230	1938.0310	1	10570.00554	3542.208390	27.712400	127.820340
[82,]	1953.0164	1954.0244	1	116.41446	38.969904	3.660130	10.647137
[83,]	1958.9972	1960.0052	1	146.46839	49.010187	3.921570	12.497593
[84,]	1963.0361	1964.0441	1	219.18520	73.321530	6.013070	12.193693
[85,]	1969.9485	1970.9565	1	300.77264	100.565706	4.444440	22.627306
[86,]	1981.0615	1982.0695	1	244.80593	81.788588	13.594800	6.016167
[87,]	1982.0622	1983.0702	1	5026.37446	1679.171935	15.686300	107.047037
[88,]	1994.0448	1995.0528	1	131.82282	44.002302	3.921570	11.220583
[89,]	1995.0238	1996.0318	1	80.08364	26.729978	3.660130	7.303013
[90,]	1998.0483	1999.0563	1	86.94219	29.013158	3.398690	8.536571
[91,]	2004.0431	2005.0511	1	76.49392	25.515850	3.660130	6.971296
[92,]	2008.0804	2009.0884	1	135.36306	45.139783	3.660130	12.332836
[93,]	2039.0826	2040.0906	1	56.21698	18.706247	2.875820	6.504665
[94,]	2052.9999	2054.0079	1	68.78940	22.867354	2.352940	9.718630
[95,]	2092.1287	2093.1367	1	81.25674	26.937889	3.660130	7.359818
[96,]	2098.0541	2099.0621	1	76.68433	25.411441	3.137260	8.099884
[97,]	2105.0082	2106.0162	1	172.43704	57.102884	3.398690	16.801439
[98,]	2109.1722	2110.1802	1	245.65955	81.315408	11.764700	6.911813
[99,]	2110.1592	2111.1672	1	5284.54564	1749.046488	14.117600	123.891206
[100,]	2126.1453	2127.1533	1	79.78905	26.365833	2.352940	11.205485
[101,]	2132.1398	2133.1478	1	59.98968	19.811183	2.352940	8.419757
[102,]	2136.1764	2137.1844	1	73.48454	24.257864	2.614380	9.278630
[103,]	2154.1284	2155.1364	1	41.38738	13.637363	2.091500	6.520374
[104,]	2157.9980	2159.0060	1	41.70303	13.735923	1.830070	7.505682
[105,]	2166.1065	2167.1145	1	73.99919	24.353337	2.352940	10.350173
[106,]	2172.1089	2173.1169	1	33.40552	10.987190	1.830070	6.003699
[107,]	2188.0666	2189.0746	1	36.97220	12.140358	1.633987	7.429897
[108,]	2211.1144	2212.1224	1	68.46039	22.420720	2.091500	10.719924
[109,]	2226.1425	2227.1505	1	102.26718	33.429187	1.633987	20.458655
[110,]	2233.0991	2234.1071	1	33.83283	11.049640	1.633987	6.762378
[111,]	2257.0072	2258.0152	1	73.28383	23.861885	1.633987	14.603468
[112,]	2283.2155	2284.2235	1	182.53729	59.238705	2.091500	28.323550
[113,]	2284.1745	2285.1825	1	24.06578	7.809114	2.091500	3.733738

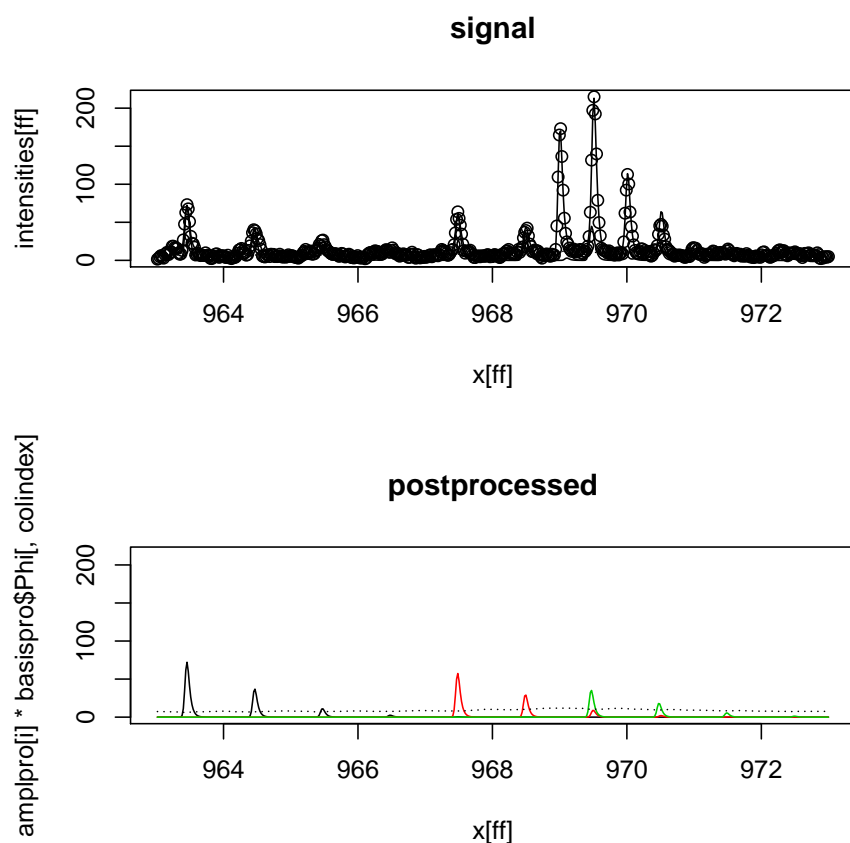
```
[114,] 2427.2620      2428.2700      1      50.99017      16.211009      1.633987      9.921134
```

```
R>
```

The results can be examined in detail graphically. We finally present some selected regions to demonstrate that our method performs well. The pre-defined method `visualize` can be used display the template fitting at several stages for regions within selected  $m/z$  intervals as specified by the arguments `lower` and `upper`.

```
R> visualize(EMGlist, x, y, lower= 963, upper = 973,
             fit = FALSE, fittedfunction = TRUE, fittedfunction.cut = TRUE,
             localnoise = TRUE, quantile = 0.5,
             cutoff.functions = 3)
```

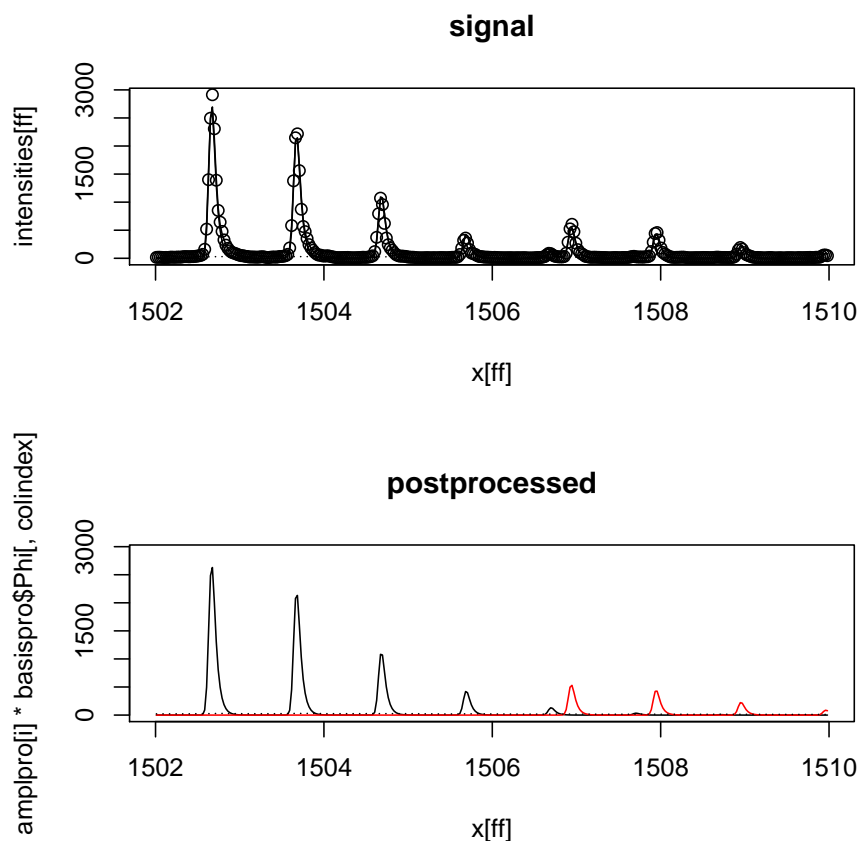
```
R>
```



```
R> visualize(EMGlist, x, y, lower= 1502, upper = 1510,
             fit = FALSE, fittedfunction = TRUE, fittedfunction.cut = TRUE,
             localnoise = TRUE, quantile = 0.5,
             cutoff.functions = 2)
```

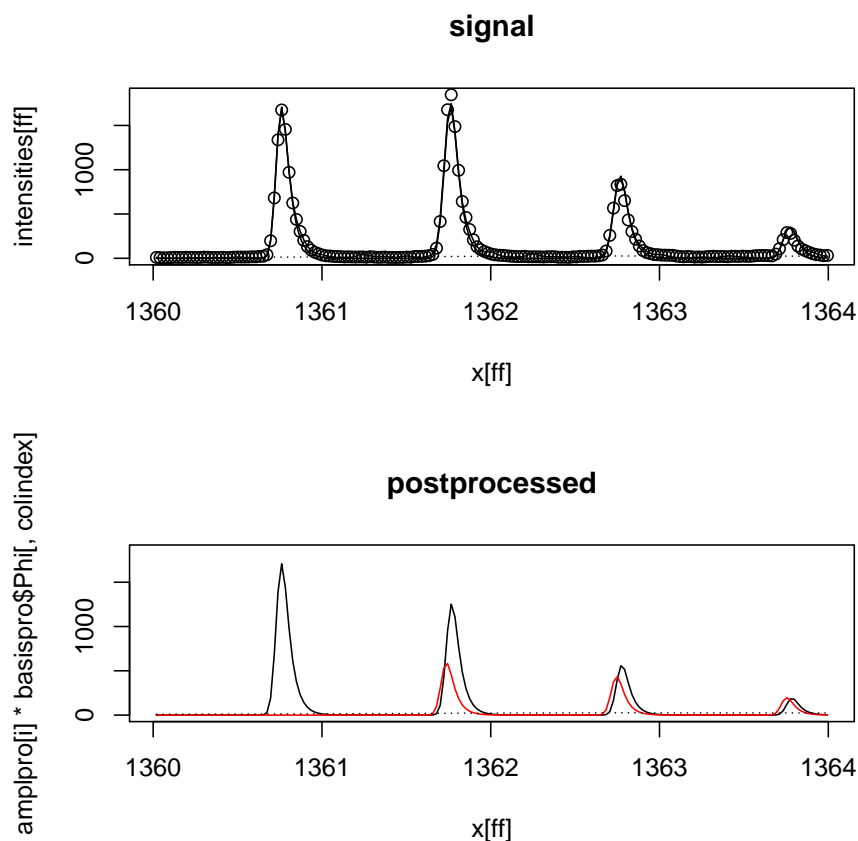
```
R>
```

```
R>
```



In the  $m/z$  range [963,973] a charge-1 peak overlaps with a more intense charge two peak. A further overlap occurs in the interval [1502,1510], and it is correctly resolved. An even more challenging problem, in which it is already difficult to unravel the overlap by visual inspection, is displayed in the following plot.

```
R> visualize(EMGlist, x, y, lower= 1360, upper = 1364,
  fit = FALSE, fittedfunction = TRUE, fittedfunction.cut = TRUE,
  localnoise = TRUE, quantile = 0.5,
  cutoff.functions = 2)
R>
```



## 4 Extension to process LC-MS runs

In the preceding sections, it has been demonstrated how IPPD can be used to process single spectrums. For LC-MS, multiple spectra, one for a sequence of retention times  $\{t_l\}_{l=1}^L$ , have to be processed. In this context, a single spectrum is referred to as scan. The resulting data can be displayed as in Figure 2 by plotting intensities over the plane defined by retention times and  $m/z$ -values. IPPD offers basic functionality to process this kind of data. Support for mzXML format as well as an implementation of the sweep line scheme as suggested in Schulz-Trieglaff et al. [2008] is provided, which is briefly demonstrated in the sequel.

```
R> directory <- system.file("data", package = "IPPD")
R> download.file("http://www.ml.uni-saarland.de/code/IPPD/CytoC_1860-2200_500-600.mzXML",
  destfile = paste(directory, "/samplefile", sep = ""),
  quiet = TRUE)
R> data <- read.mzXML(paste(directory, "/samplefile", sep = ""))
R>
```

The sweep line scheme aggregates the peaklists of multiple scans by looking for blocks of consecutive retention times at which there is signal at nearby  $m/z$ -positions. The output is a quadruple consisting of a retention time interval, a  $m/z$ -position, a charge state and

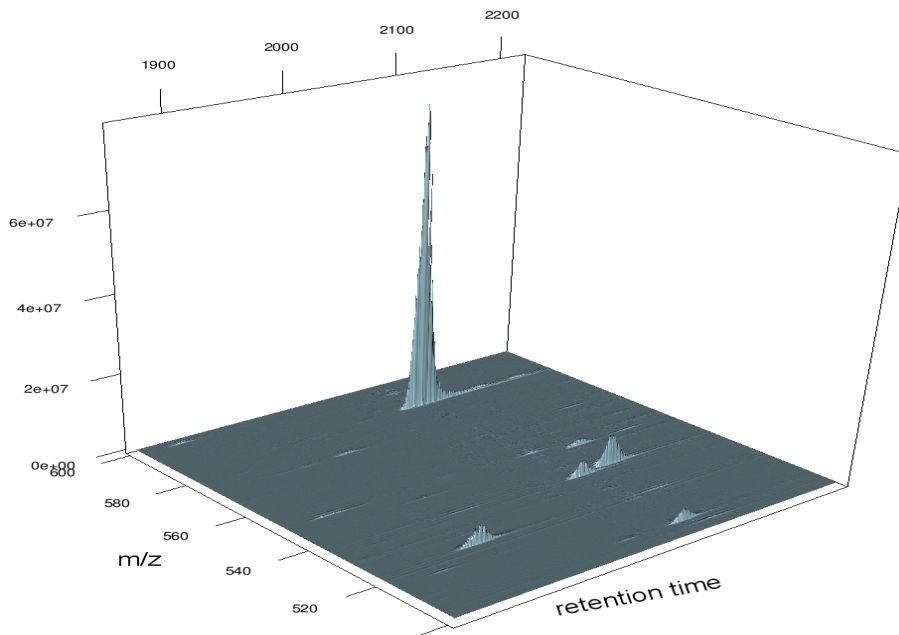


Figure 2: Graphical display of the sample `mzXML` file used in the code.

a quantification of the intensity. The intervals are found by sequentially processing the results of `getPeaklist`, where the results of each peaklist will lead to extensions of existing interval of preceding lists or to the creation of new intervals; intervals are closed once they have not been extended after processing more than `gap` additional peaklists, where `gap` is a parameter to be specified by the user. For more details, we refer to [Schulz-Trieglaff et al. \[2008\]](#). The function `analyzeLCMS` runs `getPeaklist` for each scan and then calls the function `sweepline`, which can as well be run independently from `analyzeLCMS` to aggregate the results. While there is a default setting, parameters can be changed by passing appropriate arguments.

```
R> processLCMS <- analyzeLCMS(data,
  arglist.getPeaklist = list(control.basis = list(charges = c(1,2,3))),
  arglist.threshold = list(threshold = 10),
  arglist.sweepline = list(minboxlength = 20))
R> boxes <- processLCMS$boxes
```

The output can be displayed as follows. The retention time intervals are given by the two columns `rt_begin` and `rt_end`, the corresponding `m/z`-positions are given by the column `loc`. Quantitative information is contained in the column `quant`. The output is visualized by means of a contour plot, where the contour lines depict intensities over the plane defined by `m/z`-positions and retention times. The intervals of the output are drawn as red lines.

```
R> print(boxes)
```

	loc	charge	quant	rt_begin	rt_end	npeaks	gapcount
[1,]	503.3474	1	4804929	2093.26	2113.46	26	6
[2,]	505.8353	2	58628590	2067.62	2097.79	34	13
[3,]	520.3157	1	4060317	2097.79	2118.62	22	11
[4,]	521.6749	3	174362565	1929.97	1971.22	57	6



[5,]	534.3278	2	82024006	2068.95	2099.10	35	12
[6,]	535.3383	1	242848614	2099.10	2150.15	70	9
[7,]	546.8268	2	43092586	2112.83	2137.57	37	2
[8,]	549.3673	1	10241899	1882.23	1925.06	56	8
[9,]	584.8515	2	358822314	2090.02	2155.51	91	10
[10,]	585.3650	2	17016750	2132.31	2155.51	31	5
[11,]	585.3691	2	13830263	2175.63	2193.95	25	3
[12,]	597.3784	2	10873799	1882.96	1906.53	35	1

```

R> rtlist <- lapply(data$scan, function(x)
  as.numeric(sub("([~0-9]*)([0-9|.]+)([~0-9]*)", "\\2", x$scanAttr)))
R> rt <- unlist(rtlist)
R> nscans <- length(rt)
R> npoints <- length(data$scan[[1]]$mass)
R> Y <- matrix(unlist(lapply(data$scan, function(x) x$peaks)),
  nrow = nscans,
  ncol = npoints,
  byrow = TRUE)
R> contour(rt, data$scan[[1]]$mass, Y, xlab = "t", ylab = "mz",
  levels = 10^(seq(from = 5, to = 6.75, by = 0.25)),
  drawlabels = FALSE)
R> for(i in 1:nrow(boxes))
  lines(c(boxes[i,"rt_begin"], boxes[i,"rt_end"]), rep(boxes[i,"loc"], 2), col = "red")
R>
R>

```



## References

- S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004.
- P. Du and R. Angeletti. Automatic Deconvolution of Isotope Resolved Mass Spectra Using Variable Selection and Quantized Peptide Mass Distribution. *Analytical Chemistry*, 78: 3385–3392, 2006.
- E. Grushka. Characterization of Exponentially Modified Gaussian Peaks in Chromatography. *Analytical Chemistry*, 44:1733–1738, 1972.
- V. Marco and G. Bombi. Mathematical functions for the representation of chromatographic peaks. *Journal of Chromatography A*, 931:1–30, 2001.
- B. Renard, M. Kirchner, H. Steen, J. Steen, and F. Hamprecht. NITPICK: peak identification for mass spectrometry data. *BMC Bioinformatics*, 9:355, 2008.
- O. Schulz-Trieglaff, R. Hussong, C. Gröpl, A. Hildebrandt, and K. Reinert. A Fast and Accurate Algorithm for the Quantification of Peptides from Mass Spectrometry Data. In *Proceedings of the Eleventh Annual International Conference on Research in Computational Molecular Biology (RECOMB 2007)*, pages 437–487, 2007.
- O. Schulz-Trieglaff, R. Hussong, C. Gröpl, A. Leinenbach, A. Hildebrandt, C. Huber, and K. Reinert. Computational quantification of peptides from LC-MS data. *Journal of Computational Biology*, 15:685–704, 2008.

- M. Senko, S. Beu, and F. McLafferty. Determination of Monoisotopic Masses and Ion Populations for Large Biomolecules from Resolved Isotopic Distributions. *Journal of the American Society for Mass Spectrometry*, 6:229–233, 1995.
- R. Tibshirani. Regression shrinkage and variable selection via the lasso. *Journal of the Royal Statistical Society Series B*, 58:671–686, 1996.