

Chromatography Model v 0.1

James Klimavicz

October 22, 2024

Overview

The purpose of this package is to create realistic-looking chromatogram data, and then analyze the chromatograms to detect peaks, and calculate peak parameters for each peak. The package attempts to capture many qualities of authentic chromatograms. Key features include:

- Baselines noise, with the ability to reproduce some autoregressive behavior in the baseline signal, as expected for a time series.
- Asymmetric peak shape, with peaks being modeled by exponentially-modified gaussian functions [1]. Peaks that elute later tend to be broader and shorter, but with similar areas.
- Solvent gradients and choices of solvent affect the baseline absorbance, and result in changes in the peak retention time of each compound. The change in retention time depends on the choice of solvents, percentage of each solvent, and the flow rate.
- Compound data for over 300 compounds, including UV-Vis data, have been pre-populated to produce random samples of various compounds. Samples can be defined with lists of compound names or CAS numbers
- Ability to introduce small random variations in peak elution times and peak heights, both constant across an injection and also within an injection to produce realist variations between chromatograms for otherwise identical samples and methods.

Contents

1	Systems	2
1.1	Columns	3
2	Methods	3
2.1	Instrument Methods	3
2.2	Processing Methods	4
3	Compounds	4
3.1	UV Spectra	4
3.2	Compound Library	4
3.3	Solvents	5
3.4	Solvent Library	5
3.5	Calculation of Retention Times	5
4	Injection	6
4.1	Chromatogram Creation Steps	6
4.2	Peak Detection and Quantification Steps	7

5 Chromatogram	7
5.1 Baseline	7
5.2 Peak Creation	7
5.3 Full Chromatogram	8
6 Analysis	8
6.1 Adaptive Signal Smoothing	8
6.2 Baseline Detection and Smoothing	9
6.3 Detection Algorithm	9
6.4 Peak Parameters	11
7 Samples	15
7.1 Sample Creation	15
7.2 Samples JSON	15
References	15
Appendices	15
A Future Implementation	15
B Classes	16
B.1 Injection	16
B.2 System	16
B.3 Column	16
B.4 InstrumentMethod	16
B.5 ProcessingMethod	16
B.6 Compound	16
B.7 Sample	16
B.8 Chromatogram	16
B.9 PeakFinder	16
B.10 PeakList	16
B.11 Peak	16
C Example Input JSON	16
C.1 System	16
C.2 Column	17
C.3 Instrument Method	17
C.4 Processing Method	18
D Adding to the Compound Library	19

1 Systems

Currently, a system is implemented to:

- store a system name,
- allow for a system-specific deviation in retention time δt_{sys} , e.g. those due to differences in tubing lengths between column components,
- facilitate injections onto the column, and
- allow for placing a new column with a different serial number into the system.

The **System** also stores any ****kwargs** provided to the class. Nothing is currently implemented for this, but this may be useful to have as a reference for, e.g., an output file.

1.1 Columns

The column itself is modeled by the `Column` class, which keep track of injections, n_{inj} . Upon instantiation, the column is given attributes for column-specific differences for peak asymmetry δA_{def} , broadening δB_{def} , and retention time δt_{def} .

It also contains an attribute for when a column may begin to fail, n_{fail} . After this number of injections is reached, each injection carries a probability of triggering a failure, equal to $0.00015 \cdot (n_{\text{inj}} - n_{\text{fail}})$. After failing, each injection causes slight increases in asymmetry δA_{fail} , broadening δB_{fail} , and retention times δt_{fail} , using the following formulas:

$$\begin{aligned} \rho &\in \text{Uniform}(0.95, 1.05) \\ d &= n_{\text{inj}} - n_{\text{fail}} + 1 \\ \xi &= 0.5\rho \cdot \ln \left(1 + 10^{-7} \cdot \left(\frac{x}{2} - 1 \right)^3 \right) \\ \delta A_{\text{fail}} &= 20 \cdot \xi \cdot (A - 1) \end{aligned} \tag{1}$$

$$\delta B_{\text{fail}} = \frac{\xi \cdot W}{20} \tag{2}$$

$$\delta t_{\text{fail}} = 1 + 0.01\xi \tag{3}$$

When it then comes to a specific injection, the column parameters $\delta A_{\text{col}} = \delta A_{\text{def}} + \delta A_{\text{fail}}$, $\delta B_{\text{col}} = \delta B_{\text{def}} + \delta B_{\text{fail}}$, and $\delta t_{\text{col}} = \delta t_{\text{def}} + \delta t_{\text{fail}}$ are surfaced to adjust the appearance of the chromatogram.

2 Methods

There are two types of methods: instrument methods, which determine how data is created (mimicking how data is collected in typical instruments), and processing methods, which determine how the chromatogram signal is processed.

2.1 Instrument Methods

Within the instrument methods are import parameters such as:

- solvents used
- solvent gradients and flow gradients used
- detection parameters, such as UV-Vis wavelengths.

2.1.1 Solvents

The current implementation includes the solvents water, THF, acetonitrile, methanol, and isopropanol. The solvents can be loaded into the method in any order. Typically, solvent pumps in HPLC systems handle either 2 or 4 solvents; four solvents are supported by the implementation. Additional information on solvents can be found in section 3.3.

Future implementations may allow for additional choices, such as those containing buffer, with the potential to support pH variations.

2.1.2 Solvent Gradients

The solvent gradient is implemented to allow for linear combinations of solvents are each stage. Each time step also includes a flow rate, with units of ml/min.

A future implementation will support a curve parameter, which will permit non-linear interpolation between points.

2.2 Processing Methods

The processing method includes:

- parameters that determine the peak smoothing and peak detection
- parameters that affect peak filtering (e.g. minimum height and area for peaks)
- retention time windows for identifying peaks
- calibration data for amount calculations

3 Compounds

The **Compound** class holds all compound-specific information, including compound metadata (e.g. name, CAS number, SMILES), and predicted compound properties, such as number of hydrogen bond donors and acceptors, polar surface area, molecular weight, etc. Also included are a default column volume retention, which was created to allow for slight retention time differences in isomeric compounds that have similar or identical calculated logP values. The compound class also stores concentration information from a **Sample**, and contains a **UVSpectrum** class, which allows for retrieval of a molar attenuation coefficient (absorptivity) ϵ at a provided wavelength. This value is then multiplied by the compound concentration to give a value proportional to the absorbance A .¹

3.1 UV Spectra

The **UVSpectrum** class creates For most compounds, UV-Vis spectra in the form of .jdx files were retrieved from the online NIST database [2]. Several compounds have their own UV-Vis spectra that are simply slight modifications of an existing spectrum for the sake of having the compound in the module library. The goal in using real UV spectra is to allow for more realistic chromatogram behavior, especially when multiple chromatograms or 3D spectra are used.

To create a UV-Vis spectrum that provides reasonable values of $\log \epsilon$ in the range of 200 to 800 nm, the following steps are taken:

1. Raw wavelength values and corresponding $\log \epsilon$ data is read from the .jdx file.
2. If the minimum wavelength in the .jdx file does not go down to 190 nm, extrapolate $\log \epsilon$ by simply quadratically increasing $\log \epsilon$.
3. Quadratically decrease $\log \epsilon$ for five nanometers of wavelength after the maximum wavelength in the .jdx file
4. Fit a cubic spline to the wavelength/ $\log \epsilon$ data

When request for a molar attenuation coefficient at a specific wavelength is made to the **UVSpectrum** class, $\log \epsilon$ is interpolated and the value $\epsilon = 10^{\log \epsilon}$ is returned.

3.2 Compound Library

The **CompoundLibrary** class acts as the central repository for all compound data. The library is created from a file named `compounds.csv`, which contains over 300 rows of compound data. Upon instantiation, the **CompoundLibrary** goes through this file and adds each compound to its compound list as a **Compound** object, with the **UVSpectrum** being created from the .jdx file in the `spectra` directory. Upon creation, the **CompoundLibrary** class is pickled and stored in a `cache` directory so that the library does not need to be recreated every time the module is used.

The compound library provides several methods important for sample creation, including being able to retrieve compounds by CAS number or names, and to generate random collections of compounds,

¹We note that this is not a true absorbance, since we do not have a true calculation of concentration at the peak retention time, as this is a function of peak shape, as well tubing and injection volumes.

possibly excluding a finite list of compounds already in a sample. This latter method is important if the user needs to general a sample with random assortments of compounds.

Compounds from the **CompoundLibrary** are returned as a deep copy, so that fields like the concentration and retention time may be set independently for each sample, which is crucial when running multiple samples in one go. The **UVSpectrum** is *not* deep copied, and is instead a reference to the original **UVSpectrum**.

3.3 Solvents

The **Solvent** class inherits from the **Compound** class, and has some additional properties that are crucial for determining solvent retention time.

3.4 Solvent Library

The **SolventLibrary** class inherits from the **CompoundLibrary** class. There are no additional methods currently implemented, but this library contains only **Solvent** objects.

3.5 Calculation of Retention Times

The calculation of retention time depends on many factors, including:

- compound properties,
- solvent composition and properties,
- solvent flow rate, and
- column properties (not yet implemented).

Solvent attributes included in a **Solvent** class object include hydrogen bond acidity A , hydrogen bond basicity B , solvent polarity P , solvent dipolarity π , and solvent dielectric ϵ .

A given **Compound** has attributes for a default retention column volume v_0 , an estimated partition coefficient $\log P$, total polar surface area σ , molecular weight MW , an estimated water solubility coefficient $\log S$, and number of hydrogen bond acceptors and donors, H_A and H_D . These properties were estimated using SwissADME [3].

Because the solvent profile changes over time, we note that the composite solvent attributes are functions of time. We will assume that these properties are linearly additive with respect to their component proportions. A method also specifies a flow $\eta(t)$, which may or may not remain constant over time.

Let $f(t)$ be the composite compound solvent interaction at time t . A later method will implement a more scientific model, but for now, the chosen parameters were chosen rationally (albeit not scientifically) to produce somewhat realistic changes in retention times based on solvent, including potential changes in elution order for some combinations of compounds. We first define the coefficients:

$$\begin{aligned}\alpha &= \frac{600}{MW \cdot (\sqrt{H_D} - 1)} \\ \beta &= \frac{600}{MW \cdot (\sqrt{H_A} - 1)} \\ \gamma &= 2 \cdot (3 - \log P) + \frac{\sigma}{MW} \\ \delta &= \frac{1 - \log S}{5}\end{aligned}$$

We then define $f(t)$ to be (again, based on some rational decision but mainly based on heuristic results):

$$f(t) = \frac{1}{120} (\alpha \cdot A(t) + \beta \cdot B(t) + \gamma \cdot P(t) + \delta \cdot \varepsilon(t) - 0.05 * (3 * \alpha + \beta + \gamma + \delta)) \quad (4)$$

Because elution time is dependent on the cumulative effects of solvent flow, we must account for both the solvent flow and $f(t)$. First, we have the eluted volume as a function of time, $V(t)$:

$$V(t) = \int_0^t \eta(\tau) d\tau, \quad (5)$$

so

$$dV = \eta(t) dt \quad (6)$$

The cumulative effect of solvent over time is then given by

$$\int_0^{V(t_R)} f(t) dV = \int_0^{v_R} f(t) \cdot \eta(t) dt \quad (7)$$

The actual retention volume v_R can then be found by solving the equation

$$v_R = v_0 - \int_0^{v_R} f(t) \cdot \eta(t) dt \quad (8)$$

In some cases for highly polar species, the heuristically calculated v_R may be less than one, which is non-nonsensical since it implies the compound will elute before a single column volume has flowed through the column. We therefore set v' to be the maximum of the calculated

We then convert this to an elution time t' using the column volume v_{col} and the flow rate:

$$t_R = t(v_R)/v_{col}, \quad (9)$$

where $t(v)$ is the time at which v cumulative volume has been eluted.

4 Injection

The **Injection** class is the core unit of the module. It handles the creation of the chromatogram from a sample and instrument method, and also allows for processing of peaks with a processing method. Raw chromatograms may be retrieved, as well as a JSON format of the entire injection.

4.1 Chromatogram Creation Steps

When an Injection is created, the following steps happen in order:

1. The **Sequence** object is updated with information about about the injection.
2. The **System** is updated with an adding an injection and incrementing the injection count on the column.
3. The **InjectionMethod** is loaded to:
 - (a) Determine the wavelengths and channel names to use for creating chromatograms
 - (b) Produce a solvent gradient profile for the chromatogram
 - (c) Determine the chromatogram length and sample rate.
4. A baseline is created for each chromatogram based on the wavelength and solvent (section 5.1) profile.
5. The compound retention times and peak shapes are calculated based on:

- (a) Default peak widths and asymmetry specified in the user parameters
 - (b) System-specific differences in retention time δt_{sys} (see section 1)
 - (c) Column-specific differences in retention time δA_{col} , δB_{col} , δt_{col} (section 1.1)
 - (d) Solvent profile, flow rate, and column volume (section 3.5)
6. Peak heights are adjusted based on compound UV-Vis absorbances at the wavelengths specified in the instrument method and specified amounts in the sample.
 7. Adjusted peaks are then added to the baselines.

4.2 Peak Detection and Quantification Steps

When a chromatogram is chosen for peak detection, the following steps occur:

1. An adaptive Savitzky-Golay smoothing step is first performed to generate a smoothed signal and smooth second derivative, as described in section 6.1
2. A baseline is determined as described in section 6.2, and subtracted from the smoothed signal
3. The peak detection algorithm described in section 6.3 is used to find peaks
4. Peak parameters for each peak are calculated (section 6.4)

5 Chromatogram

5.1 Baseline

The chromatogram baseline is modeled as

$$y(t) = b(t) + \mathcal{X}_t, \quad (10)$$

where $b(t)$ is a background signal dependent on the solvent profile as outlined in section 2, and \mathcal{X}_t is a one-parameter autoregressive process $\mathcal{X}_t = \varphi \mathcal{X}_{t-1} + \varepsilon_t$ with user-provided parameter φ and white noise process $\varepsilon \sim \mathcal{N}(0, \sigma_\varepsilon^2)$.

5.2 Peak Creation

This module creates peaks in the shape of an exponentially-modified gaussian (EMG), which has been used previously to model chromatography peaks [1, 4]. The module uses the `scipy.stats.exponnorm` module to calculate the EMG distribution, which uses the formula

$$f(x, K) = \frac{1}{2K} \exp\left(\frac{1}{2K^2} - \frac{x}{K}\right) \operatorname{erfc}\left(-\frac{x - \frac{1}{K}}{\sqrt{2}}\right), \quad (11)$$

where

$$\operatorname{erfc}(z) = 1 - \frac{2}{\sqrt{\pi}} \int_0^z \exp(-t^2) dt.$$

A **PeakCreator** module is created when an **Injection** object is instantiated. Upon creation, several key injection-specific attributes are created:

- injection retention time shift $\delta t_{\text{inj}} = \delta t_{\text{sys}} + \delta t_{\text{rand}}$, which is a combination of the system-specific shift δt_{sys} (see section 1), and a random time offset $\delta t_{\text{rand}} = \text{Uniform}(-\delta t_{\text{os}}, \delta t_{\text{os}})$, where δt_{os} is specified in the user parameters.
- injection height multiplier $\delta h_{\text{inj}} = \text{Uniform}(1 - \delta h_{\text{os}}, 1 + \delta h_{\text{os}})$, where δh_{os} is specified in the user parameters.

Algorithm 1 Adaptive Signal Smoothing

Require: \mathbf{y} , w_{\min} , w_{\max} , w_{var} , k $\mu_{\text{local}} \leftarrow \text{UNIFORM}(\mathbf{y}, w_{\text{var}})$ \triangleright Array of local means of size w_{var} $\sigma_{\text{local}}^2 \leftarrow \text{UNIFORM}((\mathbf{y} - \mu_{\text{local}})^2, w_{\text{var}})$ \triangleright Array of local variances of size w_{var} $\varsigma_{\text{bg}} \leftarrow$ root mean square of background standard deviation $\sigma_{\text{adj}}^2 \leftarrow (\sigma_{\text{local}} / (k \cdot \varsigma_{\text{bg}}))^2$ $\mathbf{w} \leftarrow w_{\max} - \sigma_{\text{adj}}^2 \cdot (w_{\max} - w_{\min})$, clipped to range $[w_{\min}, w_{\max}]$ adjust windows \mathbf{w} to only change one step at a timediscretize windows \mathbf{w} to nearest integer $w_{\min} < \mathbf{w}_i = 2k + 1 < w_{\max}$ $\tilde{\mathbf{y}} \leftarrow \mathbf{0}$ **for all** $i \in \{1, \dots, \text{len}(\mathbf{w})\}$ **do** $\tilde{\mathbf{y}}_i \leftarrow \text{SAVGOL}(\mathbf{y}_i, \mathbf{w}_i)$ **return** $\tilde{\mathbf{y}}$

Additionally, the **PeakCreator** handles:

- determination of broadening over time
- determination of asymmetry increase over time
- adjustment of peak high to ensure broadening/asymmetry changes do not result in changes to average peak area
- incorporation of **Column**-related retention time shifts and broadening and asymmetry increases (see section 1.1).

5.3 Full Chromatogram

The full chromatogram for each detection wavelength is generated by sequentially adding peaks to a baseline. For each compound in the injection, a compound-specific absorbance is calculated from the **Compound** class, and an exponentially-modified gaussian peak is created based on actual compound retention time (see section 3.5), and compound- and time-specific asymmetry and peak broadening. The peak shape calculations are performed with the **PeakCreator** class.

6 Analysis

The analysis portion of the module is intended to act completely independently from the chromatogram generation portion — that is, the only information used to analyze peaks comes from the chromatogram time and signal points, and the **InstrumentMethod** provided. The latter is described in section 2.2.

6.1 Adaptive Signal Smoothing

The signal smoothing algorithm was implemented as inspired by Chromeleon’s “Auto” smoothing algorithm, as described in [5]. Traditionally, smoothing has been performed using a Savitzky-Golay (SG) filter with a constant-width window across the whole chromatogram. However, choosing too small a window results in insufficient smoothing, especially in areas without any peaks, while large might provide adequate smoothing in flatter areas, but deform peaks by making them shorter and wider. Additionally, larger windows may introduce artifacts like baseline signal ringing into the chromatogram.

The goal of the adaptive smoothing is to use a maximally sized smoothing window in regions without peaks, and reduce the size of the smoothing window around peaks. While the technical note does not describe much in the way of implementation, algorithm 1 provides the outline of the implementation for this module.

The algorithm takes the signal \mathbf{y} and a max window size $n = 2k + 1, k \in \mathbb{Z}^+$ as input.

6.2 Baseline Detection and Smoothing

The baseline is detected using an asymmetric least squares algorithm. The algorithm currently implemented is the *Peaked Signal's Asymmetric Least Squares Algorithm* (psalsa) method [6].

The algorithm takes three parameters: $p \in (0, 1)$, which is a weight parameter; $s > 0$, which is a multiplicative factor that penalizes changes in the second derivative of the baseline, and $k \geq 1$, which controls the exponential decay of weights in peak regions.

The weights \mathbf{w} are calculated as:

$$\mathbf{w}_i = \begin{cases} p \cdot e^{\frac{-r_i}{k}} & \text{if } r_i > 0 \\ 1 - p & \text{otherwise} \end{cases} \quad (12)$$

for residuals $r_i = y_i - z_i$, where z_i is the baseline fit at point i . The cost function for psalsa is then

$$C(\mathbf{r}, \mathbf{w}, \mathbf{z}) = \sum_i \mathbf{w}_i r_i^2 + s \sum_i (\Delta^2 z_i)^2, \quad (13)$$

where Δ^2 second difference operator. Minimization of this cost function gives

$$\mathbf{z} = (\mathbf{W} + s \cdot \mathbf{D}^T \mathbf{D})^{-1} \mathbf{W} \mathbf{y}, \quad (14)$$

where \mathbf{D} is the tridiagonal second difference matrix, and $\mathbf{W} = \text{diag}(\mathbf{w})$. The implemented method perturbs \mathbf{D} such the $\mathbf{D}_{1,1} = \mathbf{D}_{n,n} = 1$.

For tractability, the chromatogram signal and times arrays are down-sampled to prevent massive matrices. The down-sampled points are then interpolated with a cubic spline. The algorithm is outlined in Algorithm 2.

6.3 Detection Algorithm

The detection algorithm is arguably the most complicated algorithm implemented in this module. It consists of:

1. Initial region detection
2. Region expansion
3. Peak refinement

Algorithm 2 Asymmetric Least Squares: psalsa

Require: $\mathbf{y}, k, s, p, \text{tol}$

$\mathbf{z} \leftarrow \text{mean}(\mathbf{y}) \cdot \mathbf{1}$	\triangleright initial baseline guess
$\mathbf{r} \leftarrow \mathbf{y} - \mathbf{z}$	\triangleright initial residuals
$\mathbf{w} \leftarrow \mathbf{1}$	\triangleright initialize all weights to 1
prev_loss $\leftarrow \infty$	
while not converged do	
$\mathbf{W} \leftarrow \text{diag}(\mathbf{w})$	\triangleright create diagonal matrix of weights
$\mathbf{z} \leftarrow (\mathbf{W} + s \cdot \mathbf{D}^T \mathbf{D})^{-1} \mathbf{W} \mathbf{y}$	\triangleright solve normal equation
$\mathbf{r} \leftarrow \mathbf{y} - \mathbf{z}$	\triangleright calculate updated residuals
loss $\leftarrow C(\mathbf{r}, \mathbf{w}, \mathbf{z})$	\triangleright calculate loss
if loss - prev_loss < tol then	\triangleright check for convergence
converged $\leftarrow \text{True}$	
prev_loss \leftarrow loss	
update \mathbf{w}	\triangleright update according to Eqn 12

Algorithm 3 Initial region detection

Require: $y, \omega, \varsigma_{bg}, w_{deriv}$ $y'' \leftarrow \text{SAVGOL}(y, w_{deriv}, \text{deriv}=2)$ \triangleright Second derivative array $\text{regions} \leftarrow []$ $i \leftarrow 1$ **while** $i < \text{len}(y'')$ **do** **if** $y''_i < \omega \cdot \varsigma_{bg}$ **then** $\text{start} \leftarrow i$ $i \leftarrow i + 1$ **while** $y''_i < \omega \cdot \varsigma_{bg}$ **do:** $i \leftarrow i + 1$ $\text{stop} = i$ **if** $\text{end} - \text{start} > 5$ **then** append $(\text{start}, \text{end})$ on regions

Algorithm 4 Region expansion

Require: $\tilde{y}, y, \varsigma_{bg}, w_{deriv}$ **Require:** start, end \triangleright start and end indices of a peak region from Alg. 3 $y' \leftarrow \text{SAVGOL}(y, w_{deriv}, \text{deriv}=1)$ **while** $\text{start} > 0 \ \&\& \ \tilde{y}_{\text{start}} > \varsigma_{bg} \ \&\& \ y'_{\text{start}} > 0$ **do** $\text{start}--$ **while** $\text{end} < \text{len}(\tilde{y}) \ \&\& \ \tilde{y}_{\text{end}} > \varsigma_{bg} \ \&\& \ y'_{\text{end}} < 0$ **do** $\text{end}++$

6.3.1 Initial Region Detection

The initial region detection is relatively simple. First, the root-mean-squared noise ς_{bg} of the second derivative of the spectrum is calculated. Then, all the contiguous regions where the second derivative is below $\omega \cdot \varsigma_{bg}$ for some user-specified $\omega > 1$ are recorded as peaks.

Some additional checking is performed to determine if the signal is above the linear limit for a detector, since peaks going above this limit may have some artifacts in the second derivative caused by saturation of the signal; it is undesirable for these peaks to be split from these artifacts in the second derivative.

The choice of ω by the user is important, because too high a value will result in few or no peaks being detected. However, too small a value may result in many false positive peaks.

Before proceeding on to region expansion, a quick filtering is performed to ensure that each region has a minimum width to prevent noise from presenting as peaks. The initial region detection process is outlined in Algorithm 3. The output from this algorithm is a list of potential peak regions.

6.3.2 Region Expansion

Once we have a set of regions, we must expand them, since the negative peaks of a second derivative are narrower than the original peak. We expand peaks until the peak raw signal drops below the noise threshold, or the first derivative changes sign.

6.3.3 Peak Refinement

Peaks are refined by iterating through the list of found peaks. The **ProcessingMethod** allows for minimum height and minimum area values to be provided. Any peak not meeting both of these requirements are removed from the list. The minimum height requirement may be specified as either a multiplier of the baseline noise, or an absolute height above the baseline.

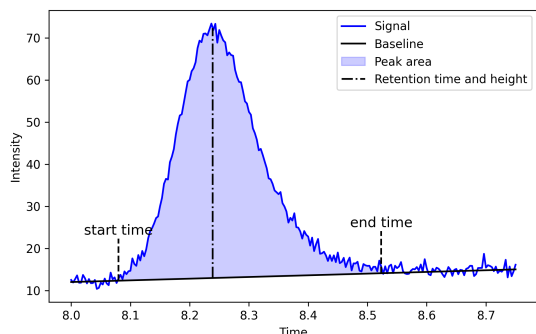


Figure 1: Peak showing peak start and end times, retention time, and baseline.

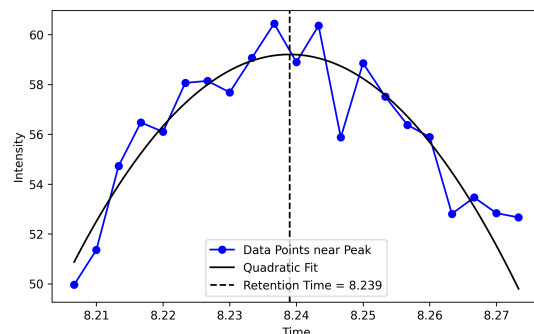


Figure 2: Peak apex showing quadratic fit at the top of the peak, and the calculated retention time.

6.4 Peak Parameters

6.4.1 Times and Signals

Peak start and end are determined by the peak detection algorithm. They are selected so the peaks may start and end at the same time, but no peaks may overlap.

start_time and **end_time**: The start and end of the peak as determined by the detection algorithm in 6.3, depicted in Figure fig:times. In the following peak parameters, start time is typically denoted t_i (for initial time) and end time as t_f (for final time).

retention_time: The time at which a peak exhibits maximum height above baseline. To account for noise in the signal, the highest point in the raw signal is found, and a quadratic curve is fit to the series of time/signal points starting three points before the maximal signal and ending three after. The time coordinate of the vertex of this quadratic is set to be the retention time, t_R . The signal coordinate of this vertex is set to be the peak height.

start_baseline, **end_baseline**: The values of the baseline at the **start_time** and **end_time**, with the units of the signal.

retention_baseline: The baseline value at the retention time, as calculated by linear interpolation from the start and end times and start and end baseline values.

start_signal and **end_signal**: The start and end signal values, with the units of the signal.

6.4.2 Type

6.4.3 Quantification

area: The area A of the peak under the raw signal and above the baseline:

$$A = \int_{t_i}^{t_f} f(t) \approx \sum_{k=u}^k f(t_k) \Delta t,$$

where $f(t)$ is the baseline-corrected signal. Units are in signal \cdot time; for UV-Vis, the units are mAU \cdot min.

relative_area: The percentage of area of the peak relative to the summed area of all peaks; that is, if there are n peaks, peak i has relative area

$$A_{\text{rel}} = 100\% \cdot \frac{A_i}{\sum_{k=1}^n A_k}$$

with units of percent.

capillary_electrophoresis_area: Equal to A/t_R , the capillary electrophoresis area is the area adjusted for migration (retention) time in capillary electrophoresis to compensate for changes in migration time between runs

height: The maximal height h of the above the baseline at the retention time t_R . This height is calculated using a quadratic approximation to the apex of the peak; see the **retention_time** entry in section 6.4.1. Height units are the units of the signal.

relative_height The percentage of the peak's height above baseline relative to the summed heights above baseline for all peaks:

$$h_{\text{rel}} = 100\% \cdot \frac{h_i}{\sum_{k=1}^n h_k}$$

with units of percent.

6.4.4 Widths

When possible, peak widths are calculated based on actual signal data, determining at what time a signal crosses a height threshold to the left and right of the retention time. For peaks that are not baseline resolved, it may be necessary to use extended peak calculations to estimate peak widths.

width_50_<left,right,full>: The left, right, and full peak width of the peak at 50% height.

width_10_<left,right,full>: The left, right, and full peak width of the peak at 10% height.

width_5_<left,right,full>: The left, right, and full peak width of the peak at 5% height.

width_4_sigma_<left,right,full>: The left, right, and full peak width of the peak at $4\sigma \approx 13.4\%$ height.

width_5_sigma_<left,right,full>: The left, right, and full peak width of the peak at $5\sigma \approx 4.4\%$ height.

width_baseline_<left,right,full>: The left and right points of inflection of the fitted exponentially-modified gaussian are first calculated. The slopes at these two points are then calculated, and the lines with this slope going through these two points are extended to the baseline. The left, right, and full baseline widths are determined based on where these slope lines intersect the baseline, as well as a vertical line at t_R .

extended_width_calculation: an array of the string value names for which extended peak width calculations were needed to calculate, as opposed to using raw signal values.

6.4.5 Statistical Moments

Statistical moments quantify aspects of the peak shape. For the calculation of peak shapes, the function $f(t)$ is the baseline-corrected signal, which is equal to the raw signal after subtraction of the baseline. For the calculation of moments, the absolute value of this function is used to avoid negative moments, which can lead to runtime errors when calculating higher moments. This will result in differences between **moment_0** and the **area** for smaller peaks when noise levels cause the signal to drop below the baseline.

moment_0: Equal to peak area. For peak start time t_i and peak end t_f and baseline-corrected signal $f(t)$,

$$\mu_0 = \int_{t_i}^{t_f} |f(t)| dt \approx \sum_{k=i}^f |f(t_k)| \Delta t.$$

Units are equal to the signal value multiplies by time in minutes, e.g. mAU · min.

moment_1: Average retention time. For peak start time t_i and peak end t_f and baseline-corrected signal $f(t)$,

$$\mu_2 = \frac{1}{\mu_0} \int_{t_i}^{t_f} t \cdot |f(t)| dt \approx \sum_{k=i}^f t_k \cdot |f(t_k)| \Delta t$$

Units are in minutes.

moment_2: Retention time variance. For peak start time t_i and peak end t_f and baseline-corrected signal $f(t)$,

$$\mu_2 = \frac{1}{\mu_0} \int_{t_i}^{t_f} (t - \mu_1)^2 \cdot |f(t)| dt \approx \sum_{k=i}^f (t_k - \mu_1)^2 \cdot |f(t_k)| \Delta t$$

Units are in min².

standard.deviation: Retention time standard deviation. Standard deviation $\sigma = \sqrt{\mu_2}$. Units are in min.

moment_3: Third central moment of retention time. For peak start time t_i and peak end t_f and baseline-corrected signal $f(t)$,

$$\mu_3 = \frac{1}{\mu_0} \int_{t_i}^{t_f} (t - \mu_1)^3 \cdot |f(t)| dt \approx \frac{1}{\mu_0} \sum_{k=i}^f (t_k - \mu_1)^3 \cdot |f(t_k)| \Delta t$$

Units are in min³.

moment_3.normalized: Third normalized moment of retention time, or statistical skewness (not to be confused with USP skewness). For peak start time t_i and peak end t_f and baseline-corrected signal $f(t)$,

$$\tilde{\mu}_3 = \frac{1}{\mu_0 \cdot \mu_2^{3/2}} \int_{t_i}^{t_f} (t - \mu_1)^3 \cdot |f(t)| dt \approx \frac{1}{\mu_0 \cdot \mu_2^{3/2}} \sum_{k=i}^f (t_k - \mu_1)^3 \cdot |f(t_k)| \Delta t$$

The normalized moment is unitless.

moment_4: Fourth central moment of retention time. For peak start time t_i and peak end t_f and baseline-corrected signal $f(t)$,

$$\mu_4 = \frac{1}{\mu_0} \int_{t_i}^{t_f} (t - \mu_1)^4 \cdot |f(t)| dt \approx \frac{1}{\mu_0} \sum_{k=i}^f (t_k - \mu_1)^4 \cdot |f(t_k)| \Delta t$$

Units are in min⁴.

moment_4.normalized: Fourth central moment of retention time, or kurtosis. For peak start time t_i and peak end t_f and baseline-corrected signal $f(t)$,

$$\tilde{\mu}_4 = \frac{1}{\mu_0 \cdot \mu_2^2} \int_{t_i}^{t_f} (t - \mu_1)^4 \cdot |f(t)| dt \approx \frac{1}{\mu_0 \cdot \mu_2^2} \sum_{k=i}^f (t_k - \mu_1)^4 \cdot |f(t_k)| \Delta t$$

The normalized moment is unitless.

6.4.6 Asymmetry

Asymmetry provides a measure for how much a peak is tailing or fronting. This metric may be used to monitor column quality. For the USP/EP and AIA definitions, as well as skewness, ideal peaks have a value near 1, while peaks less than 1 signify a fronting peak, and values greater than 1 signify a tailing peak. Unless otherwise stated in a test or assay, the USP standard requires **asymmetry_USP** to fall between 0.8 and 1.8 if the peak is to be used for quantification.

For statistical asymmetry, ideal peaks are scattered around 0; this value is negative for fronting peaks and positive for tailing peaks.

asymmetry_USP: US Pharmacopoeia calculation for asymmetry; this is identical to the European Pharmacopoeia (EP) and Japanese Pharmacopoeia (JP) definition:

$$A = \frac{RW_{5\%} + LW_{5\%}}{2 \cdot LW_{5\%}} = \frac{W_{5\%}}{2 \cdot LW_{5\%}}$$

for 5% left, right and full widths, $LW_{5\%}$, $RW_{5\%}$, and $W_{5\%}$. This value is undefined if these widths cannot be calculated. This value is also sometimes called the **tailing factor** or **asymmetry factor**.

asymmetry_AIA: AIA calculation for asymmetry:

$$A = \frac{RW_{10\%}}{LW_{10\%}}$$

for 10% left and right widths, $LW_{10\%}$ and $RW_{10\%}$. This value is undefined if these widths cannot be calculated.

asymmetry_statistical: Calculation of asymmetry using moments:

$$A = \frac{\mu_1 - t_R}{\sqrt{\mu_2}}$$

for first and second statistical moments μ_1 and μ_2 (see 6.4.5), and retention time t_R .

skewness: Metric similar to asymmetry. Not to be confused with statistical skewness, which is equal to the third normalized moment of retention time, $\tilde{\mu}_3$, as defined in 6.4.5. This measure of skewness is defined as:

$$S = \frac{RW_{10\%} + LW_{10\%}}{2 \cdot LW_{10\%}}$$

for 10% left and right widths, $LW_{10\%}$ and $RW_{10\%}$. This value is undefined if these widths cannot be calculated.

6.4.7 Resolution

Resolution measures how well two peaks are separated, using the width of two peaks and the distance in time between them. Resolution is always defined with regard to a reference peak; a peak cannot reference itself. The module uses the USP, EP/JP, and statistical moments methods each for up to five different reference peaks, where the string value in parenthesis replaces **<type>** in the resolution field name:

1. previous main peak (**previous_main**): the previous quantified peak
2. next main peak (**next_main**): the next quantified peak
3. previous named peak (**previous_named**): the previous peak with a name that is not equal to **None**, **"**, or **"unknown"**.
4. next named peak (**next_named**): the next peak with a name that is not equal to **None**, **"**, or **"unknown"**.
5. specified reference peak (**reference**): a peak specified by name in the processing method (see section 2.2).

resolution_ep<type>: European Pharmacopoeia calculation for resolution R , where for a reference peak ref ,

$$R = 1.18 \left| \frac{t_{ref} - t_R}{W_{50\%,ref} + W_{50\%,R}} \right|$$

resolution_usp<type>: US Pharmacopoeia calculation for resolution R , where for a reference peak ref ,

$$R = 2 \left| \frac{t_{ref} - t_R}{BW_{ref} + BW_R} \right|$$

resolution_statistical<type>: Statistical moments calculation for resolution R , where for a reference peak ref ,

$$R = \left| \frac{t_{ref} - t_R}{2 (\sqrt{\mu_{2,ref}} + \sqrt{\mu_2})} \right|$$

signal_to_noise:

6.4.8 Theoretical Plates

plates_EP: European Pharmacopoeia calculation for theoretical plates TP, where for retention time t_R and 50% width $W_{50\%}$,

$$TP = 5.54 \left(\frac{t_R}{W_{50\%}} \right)^2$$

plates_USP: US Pharmacopoeia calculation for theoretical plates TP, where for retention time t_R and baseline width BW,

$$TP = 16 \left(\frac{t_R}{BW} \right)^2$$

plates_statistical: Statistical Moments calculation for theoretical plates TP, where for retention time t_R and second statistical moment μ_2 ,

$$TP = \frac{t_R^2}{\mu_2}$$

7 Samples

7.1 Sample Creation

7.1.1 Product Stability

7.1.2 Reaction Optimization

7.2 Samples JSON

References

- [1] Pamela J Naish and S Hartwell. “Exponentially modified Gaussian functions—a good model for chromatographic peaks in isocratic HPLC?” In: *Chromatographia* 26.1 (1988), pp. 285–296. DOI: 10.1007/BF02268168.
- [2] Victor Talrose et al. “UV/Visible Spectra”. In: *NIST Chemistry WebBook, NIST Standard Reference Database Number 69*. Ed. by P.J. Linstrom and W.G. Mallard. Gaithersburg MD, 20899: National Institute of Standards and Technology. DOI: 10.18434/T4D303.
- [3] A. Daina, O. Michielin, and V. Zoete. “SwissADME: a free web tool to evaluate pharmacokinetics, drug-likeness and medicinal chemistry friendliness of small molecules”. In: *Sci Rep* 7 (2017), p. 42717. DOI: 10.1038/srep42717.
- [4] Yuri Kalambet et al. “Reconstruction of chromatographic peaks using the exponentially modified Gaussian function”. In: *Journal of Chemometrics* 25.7 (2011), pp. 352–356. DOI: 10.1002/cem.1343.
- [5] Peter Zipfell and Darren Barrington-Light. *Technical Note 70698: Intelligent integration using Cobra and SmartPeaks*. Tech. rep. Thermo Fisher Scientific Inc., 2019. URL: <https://assets.thermofisher.com/TFS-Assets/CMD/Technical-Notes/TN-70698-CDS-SmartPeaks-Cobra-TN70698-EN.pdf>.
- [6] Sergio Oller-Moreno et al. “Adaptive Asymmetric Least Squares baseline estimation for analytical instruments”. In: *2014 IEEE 11th International Multi-Conference on Systems, Signals & Devices (SSD14)*. 2014, pp. 1–5. DOI: 10.1109/SSD.2014.6808837.

Appendices

A Future Implementation

Future iterations of this module will include (in no particular order):

- A method to allow for different column types to produce different elution profiles
- Enablement of pH-dependent phenomena, e.g. peak broadening and change in retention time, based on buffer pH and compound pK_a and pK_b values.
- Improved scientific model for estimating elution times based on solvent, column, and compound properties
- Additional spectra types
 - 3D UV-Vis
 - CAD/ELSD
 - Fluorescence
 - MS

B Classes

B.1 Injection

B.1.1 Instantiation

```

1 Injection(
2     sample: Sample,
3     method: InstrumentMethod,
4     processing_method: ProcessingMethod,
5     sequence: Sequence,
6     system: System,
7     user: str | None = "admin",
8     injection_time: datetime.datetime | None = None,
9 )

```

B.1.2 Methods

B.1.3 Dictionary Keywords

B.2 System

B.3 Column

B.4 InstrumentMethod

B.5 ProcessingMethod

B.5.1 Instantiation

B.5.2 Methods

B.5.3 Dictionary Keywords

B.6 Compound

B.7 Sample

B.8 Chromatogram

B.9 PeakFinder

B.10 PeakList

B.11 Peak

C Example Input JSON

C.1 System


```

1 {
2   "name": "Curie",
3   "software": "Chromeleon",
4   "manufacturer": "ThermoFisher",
5   "modules": [
6     {
7       "name": "Sampler",
8       "manufacturer": "Thermo Fisher Scientific Inc.",
9       "type": "Autosampler",
10      "detector_type": null,
11      "part_number": "ACC-3000T",
12      "serial_number": "23056882"
13    },
14    {
15      "name": "PumpModule.Pump",
16      "manufacturer": "Thermo Fisher Scientific Inc.",
17      "type": "Pump",
18      "detector_type": null,
19      "part_number": "LPG-3400SD",
20      "serial_number": "96014241"
21    },
22    {
23      "name": "UV",
24      "manufacturer": "Thermo Fisher Scientific Inc.",
25      "type": "Detector",
26      "detector_type": "UV-Vis",
27      "part_number": "DAD-3000RS",
28      "serial_number": "54202953",
29      "firmware_version": "1.10.0",
30      "driver_version": "7.3.2.10759"
31    }
32  ],
33  "column": {}
34 },
35 "system-retention-time-offset": 0.021
36 }

```

C.2 Column

```

1 {
2   "inner_diameter": {
3     "value": 10,
4     "unit": "mm"
5   },
6   "length": {
7     "value": 150,
8     "unit": "mm"
9   },
10  "type": "C18",
11  "serial_number": "1995032",
12  "injection_count": 0,
13  "particle_size": {
14    "value": 5,
15    "unit": "um"
16  }
17 }

```

C.3 Instrument Method

```

1 {
2   "name": "column-quality-check",
3   "run_time": 8,
4   "mobile-phases": [
5     {
6       "name": "water",
7       "id": "A"
8     },
9     {
10      "name": "THF",
11      "id": "B"
12    },
13    {
14      "name": "methanol",
15      "id": "C"
16    },
17    {
18      "name": "acetonitrile",
19      "id": "D"
20    }
21  ],
22  "mobile-phase-gradient-steps": [
23    {
24      "time": 0.0,
25      "flow": 5.0,
26      "percent_a": 100.0,
27      "percent_b": 0.0,
28      "percent_c": 0.0,
29      "percent_d": 0.0,
30      "curve": 5
31    },
32    {
33      "time": 8.0,
34      "flow": 5.0,
35      "percent_a": 85.0,
36      "percent_b": 15.0,
37      "percent_c": 0.0,
38      "percent_d": 0.0,
39      "curve": 5
40    }
41  ]
42 }

```

```

41   ],
42   "detection": {
43     "uv_vis_parameters": [
44       {
45         "name": "UV-VIS_1",
46         "wavelength": 250,
47         "bandwidth": 4
48       },
49       {
50         "name": "UV-VIS_2",
51         "wavelength": 270,
52         "bandwidth": 4
53       }
54     ]
55   },
56   "creation": {
57     "computer": "CHROMELEON2",
58     "comment": null,
59     "time": "2020-08-31T12:18:06.01Z",
60     "user": "cmadmin",
61     "data_vault_name": "CHROMELEON2/ChromeleonLocal"
62   },
63   "last_update": {
64     "computer": "CHROMELEON2",
65     "comment": null,
66     "time": "2020-08-31T12:24:33.75Z",
67     "user": "cmadmin",
68     "data_vault_name": "CHROMELEON2/ChromeleonLocal"
69   },
70   "sample_introduction": {
71     "dilution_factor": 1.0,
72     "injection_volume": 20.0,
73     "weight": 1.0,
74     "internal_standard_amount": 1.0,
75     "autodilution_ratio": 0.0
76   }
77 }

```

C.4 Processing Method

```

1  {
2    "name": "column performance test",
3    "detection_parameters": {
4      "background_noise_range": {
5        "minimum": 0,
6        "maximum": 250
7      },
8      "noise_threshold_multiplier": 10,
9      "peak_limit": 5,
10     "minimum_height": {
11       "type": "baseline_noise_multiplier",
12       "value": 10
13     },
14     "minimum_area": 0.3
15   },
16   "resolution_reference": "TS-8391",
17   "peak_identification": [
18     {
19       "min_time": 2.0,
20       "max_time": 2.5,
21       "method": "largest",
22       "name": "TS-8391"
23     },
24     {
25       "min_time": 3.6,
26       "max_time": 4.1,
27       "method": "largest",
28       "name": "klimavizinib",
29       "calibration": {
30         "type": "linear",
31         "amount_unit": "umol/mL",
32         "points": [
33           {
34             "area": 101.2,
35             "amount": 2.5
36           },
37           {
38             "area": 99.2,
39             "amount": 2.5
40           },
41           {
42             "area": 1002,
43             "amount": 25
44           },
45           {
46             "area": 987.4,
47             "amount": 25
48           },
49           {
50             "area": 10,
51             "amount": 0.25
52           },
53           {
54             "area": 9.3,
55             "amount": 0.25
56           },
57           {
58             "area": 9.2,
59             "amount": 0.25
60           },
61           {
62             "area": 0.913,
63             "amount": 0.025

```

```

64         },
65         {
66             "area": 1.01,
67             "amount": 0.025
68         },
69         {
70             "area": 4.989,
71             "amount": 0.125
72         },
73         {
74             "area": 5.12,
75             "amount": 0.125
76         }
77     ]
78 },
79 {
80     "min_time": 5.9,
81     "max_time": 6.4,
82     "method": "largest",
83     "name": "caffeine"
84 },
85 ]
86 }
87 }

```

D Adding to the Compound Library

The compound library is stored in the