

ООО Амперсенд

КОМПЬЮТЕРНАЯ АВТОМАТИЗАЦИЯ ХРОМАТОГРАФИИ



[www.mullichrom.ru](http://www.mullichrom.ru)  
[support@ampersand.ru](mailto:support@ampersand.ru)  
+7(499)196-52-90

## Сколько нужно точек, чтобы правильно измерить площадь пика?

Юрий Каламбет, ООО «Амперсенд», Москва

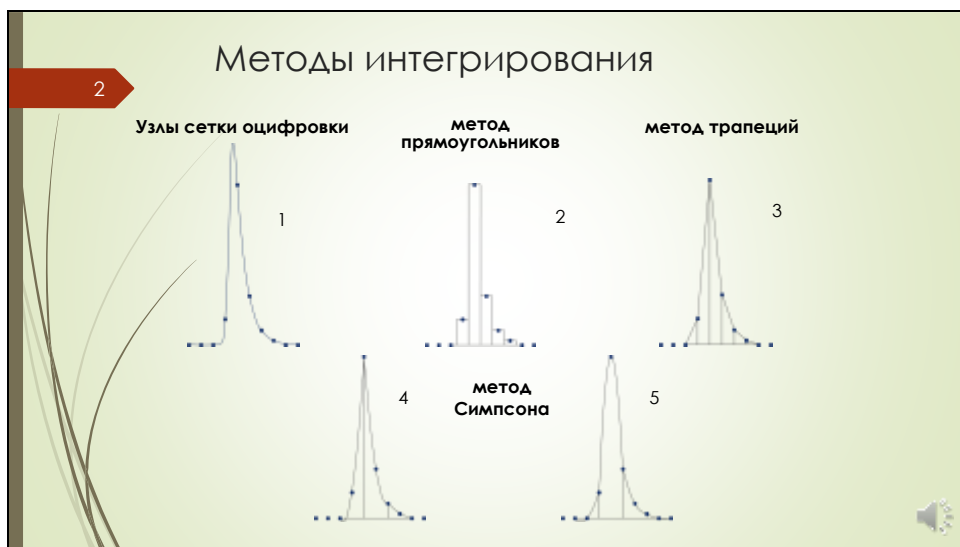
Юрий Козьмин, Институт биоорганической химии, Москва

Андрей Самохин, Московский Государственный Университет, Химический факультет

[Kalambet@ampersand.ru](mailto:Kalambet@ampersand.ru)



Здравствуйте! Меня зовут Юрий Каламбет, и сегодня расскажу вам о том, сколько точек нужно для правильного измерения площади пика. Этот вопрос касается далеко не только хроматографии, но разберем мы его на примере хроматографических пиков, поскольку авторы публикации работают в области хроматографии.



Начнем с некоторых понятий теории измерения. Обычно сигнал измеряется (переводится в цифровую форму) через равные промежутки времени. На рисунке 1 показана модель характерной формы хроматографического сигнала с отмеченными координатами оцифрованных точек. По ординате отложен отклик детектора, по абсциссе – время. Одна из основных задач обработки данных – оценить площадь под кривой по координатам измеренных точек, процесс оценки называется интегрированием.

На рисунке 2 показано интегрирование методом прямоугольников. Каждая оцифрованная точка задает прямоугольник, площадь равна сумме площадей прямоугольников. На рисунке 3 проиллюстрировано интегрирование методом трапеций. В этом методе все оцифрованные точки соединяются ломаной линией и оценкой площади считается площадь под ломаной. В методе Симпсона, проиллюстрированном на рисунках 4 и 5, каждые три соседние точки соединяются параболой, и площадь пика оценивается как площадь под соединенными параболой. Обратим внимание, что в методе Симпсона можно построить две разных оценки площади, начав конструирование парабол с четной или нечетной точки.

## Большинство методов для отдельно стоящего пика даёт один ответ = **правило прямоугольника**

Все методы интегрирования можно записать одной формулой

$$A = \Delta x \cdot \sum w(x_i) \cdot f(x_i)$$

и отличаются они друг от друга только весовыми коэффициентами  $w(x_i)$ . Ниже приведены наборы весов точек, получающиеся при интегрировании по разным правилам

- Правило прямоугольника: [001111...111100]
- Правило трапеций: [001222...222100]/2
- Правило Симпсона, усредненное по двум вариантам начала интегрирования:  
 ([014242...424100]/3+  
 [001424...242410]/3)/2=  
 [015666...666510]/6
- Разница между правилами интегрирования сводится к разному учету периферийных точек, значения в которых для отдельно стоящего пика равны нулю
- Большинство правил дают один и тот же результат: **площадь по правилу прямоугольника**

Площадь пика вычисляется по измерениям как сумма их ординат, взятая с некоторыми весами. Разные методы интегрирования отличаются набором весов. Аккуратные расчеты показывают, что три рассмотренных метода – прямоугольников, трапеций и Симпсона – для отдельно стоящих пиков дадут одинаковый результат, поскольку все отличия сконцентрированы на краях пика, где функция близка к нулю.

## Откуда берутся узкие пики в хроматографии?

- Медленное измерение: «быстрые» сканирующие детекторы, ГХ-МС, ЖХ-МС
- Быстрая хроматография: например, второе D в 2D-хроматографии, жидкостная хроматография сверхвысокого давления (UPLC)

Собственно, откуда берутся узкие пики в хроматографии? И что значит узкие? Типовая задача в аналитической химии – измерение площади пика для оценки концентрации вещества в растворе. Число измерений в единицу времени определяется природой детектора. Какое бы ни использовалось оборудование, у него есть свой предел, так что при ширине пика, сравнимой с постоянной времени измерения, результаты, полученные с использованием прибора, становятся недостоверными. В качестве примера можно привести газовую хроматографию с квадрупольным масс-спектрометрическим детектором. Этот детектор делает развертку масс по времени, и чем шире диапазон масс, тем реже происходит измерение каждой индивидуальной массы. Задача повышения информативности (диапазон масс) вступает в конфликт с задачей точного измерения площади.

## Оценка погрешности интегрирования согласно работе Нормана Дайсона (N.Dyson, J.Chromatography A, 1999, 842, 321-340)

- Погрешность интегрирования методом трапеций может быть оценена по формуле

$$I_{true} - I_{meas} = \left( \frac{W_b^3}{12n^2} \right) |h''(t)| = (\epsilon^3 N / 12) \cdot |f''(t)| \quad (1)$$

- Погрешность интегрирования по методу Симпсона меньше, чем по методу прямоугольников или трапеций
- Для надежного измерения площади требуется от 40 до 100 точек на пик
- Асимметричные пики требуют гораздо большего числа точек, чем пики Гауссовой формы при одинаковой погрешности интегрирования

Норман Дайсон в середине 90-х годов написал очень полезную и популярную книгу «Методы интегрирования в хроматографии». Эта книга выдержала несколько переизданий. К сожалению, он также написал статью, посвященную интегрированию «очень узких» хроматографических пиков, некоторые утверждения из которой приведены на слайде. Все утверждения, кроме последнего, неверны. Самый главный и неприятный вывод, имеющийся в этой статье, состоит в том, что для правильного измерения площади нужно иметь очень высокую частоту оцифровки, не ниже 40 точек на пик. Есть работы других авторов с отличающимися оценками, но четкого вывода оценки необходимой частоты измерения мы не нашли.

## «Пикоподобная» функция $f(x)$

- Непрерывно и бесконечно дифференцируемая
- Функция и все ее производные равны нулю вне пределов ограниченного интервала значений аргумента  $x$
- Свойство: Определенный интеграл любой производной по всей области пика равен нулю. (Интеграл производной равен производной предыдущего порядка; все производные по краям пика обращаются в ноль.)
- Хроматографический пик – «пикоподобная» функция
- Пик можно рассматривать как профиль распределения вероятности выхода молекул вещества из колонки. Мера его ширины – стандартное отклонение распределения, корень из дисперсии распределения

Не буду вдаваться в детали математических вычислений – желающие могут приостановить демонстрацию и почитать текст слайдов подольше. Обратим внимание только на то, что частоту измерений мы сравниваем не с шириной пика, поскольку у нас нет устоявшегося удобного определения ширины пика, а со стандартным отклонением соответствующего пику распределения вероятностей, которое обозначается греческой буквой сигма. В зависимости от определения, ширина пика составляет от 4 до 8 сигма.

## Оценка погрешности интегрирования

- Ряд Тейлора

$$f(x + \tau) = f(x_i) + f'(x_i) \cdot \tau + \frac{1}{2} f''(x_i) \cdot \tau^2 + \dots + \frac{1}{n!} f^{(n)}(x_i) \cdot \tau^n + \dots; \quad -\frac{\varepsilon}{2} < \tau < \frac{\varepsilon}{2}$$

- Для оценки площади  $A$  пика интегрируем ряд Тейлора в  $\varepsilon/2$ -окрестности узла и суммируем оценки всех узлов; координаты узлов  $x_i = \frac{\varepsilon}{2} + i \cdot \varepsilon$

$$A = A_0 + \Delta A_2 + \Delta A_4 + \dots$$

$$A_0 = \varepsilon \cdot \sum_{i=1}^N f(x_i)$$

$$\Delta A_2 = \frac{2}{3!} \left(\frac{\varepsilon}{2}\right)^3 \sum_{i=1}^N f''(x_i)$$

...

$$\Delta A_{2k} = \frac{2}{(2k+1)!} \left(\frac{\varepsilon}{2}\right)^{2k+1} \cdot \sum_{i=1}^N f^{(2k)}(x_i) = K_{2k} \cdot \sum_{i=1}^N f^{(2k)}(x_i)$$

- Попробуем оценить средний (по положению начальной точки сетки) вклад каждого члена суммы



Мы попытались вывести правильную оценку погрешности измерения площади,

## Средний (по положению начальной точки сетки) вклад $2k$ -слагаемого

- Среднее значение вклада  $2k$ -й производной

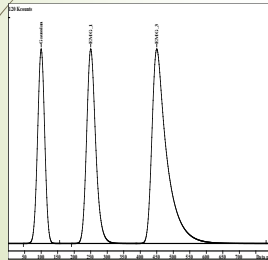
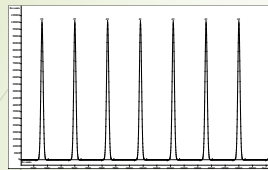
$$\begin{aligned} \Delta A_{2k} &= \frac{1}{\varepsilon} \int_{-\frac{\varepsilon}{2}}^{\frac{\varepsilon}{2}} \Delta A_{2k}(\tau) d\tau = \frac{1}{\varepsilon} K_{2k} \int_{-\frac{\varepsilon}{2}}^{\frac{\varepsilon}{2}} \sum_{i=1}^N f^{(2k)}(x_i + \tau) d\tau \\ &= \sum_{i=1}^N \int_{-\frac{\varepsilon}{2}}^{\frac{\varepsilon}{2}} f^{(2k)}(x_i + \tau) d\tau = \int_0^{(N+1)\varepsilon} f^{(2k)}(\tau) d\tau \\ &= f^{(2k-1)}(0) - f^{(2k-1)}((N+1) \cdot \varepsilon) = 0 \end{aligned}$$

- Среднее значение вклада каждой производной равно нулю
- Следовательно,  $A_0$  (площадь по методу прямоугольников) является несмещенной оценкой площади пика

И выяснили, что оценка площади по методу прямоугольников является несмещенной. Это утверждение означает, что среднее значение площади, посчитанной методом прямоугольников, по всем возможным положениям первой точки сетки оцифровки, равно истинной площади пика.



## Моделирование пиков



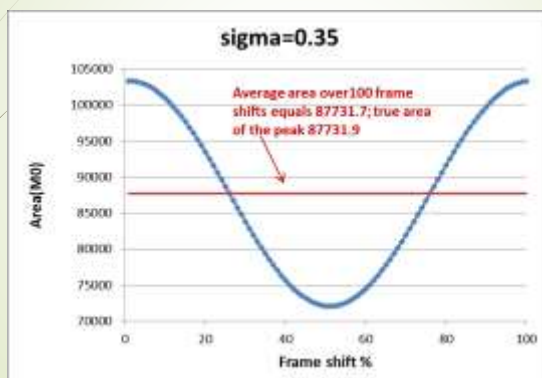
100 неперекрывающихся пиков  
 Без шума  
 Экспоненциально модифицированная  
 гауссиана  
 Высота=100000 ед.  
 Отклик в каждой точке округлен до целого числа  
 $\sigma_G=0.35\dots 8$  точек  
 $t/\sigma_G = 0$ (Гауссиана); 1(ЭМГ-1); 3(ЭМГ-3)

Шаг = Целое+0.01

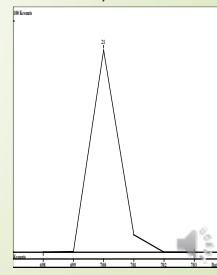
Свойство ЭМГ: Дисперсия ЭМГ  $\sigma^2=\sigma_G^2+t^2$

Теорию несложно проверить простым модельным экспериментом. Делаем 100 отдельно стоящих пиков, у каждого следующего вершина сдвинута относительно предыдущего на целое число плюс 1/100 точки, и меряем их площади. У нас должна получиться средняя площадь, равная истинной площади пика. Истинная площадь пика известна из его формулы. Процесс был повторен три раза для пиков разной формы: симметричного (Гауссиана), слабо асимметричного (ЭМГ-1) и сильно асимметричного (ЭМГ-3).

## ЗАВИСИМОСТЬ ПЛОЩАДИ ОТ СДВИГА СЕТКИ



Гауссиана;  
Площадь vs сдвиг  
сетки. Сдвиг на 100%  
эквивалентен сдвигу  
на 1 точку.



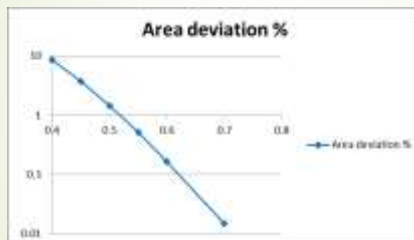
На графике приведена зависимость площади пика (по вертикали) от сдвига первой точки сетки оцифровки. Сдвиг на 100% соответствует сдвигу по времени на 1 точку.

Генерировались пики Гауссовой формы, частота оцифровки равна 0.35 точек на сигму.

Коричневая линия соответствует истинной площади пика. Форма одного из сгенерированных пиков приведена справа.

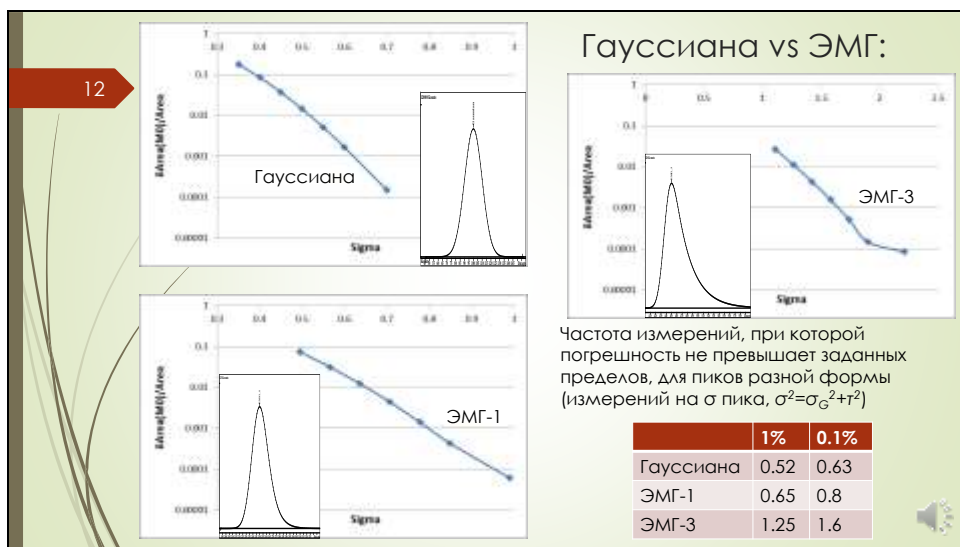
Как и ожидалось, среднее значение площади по 100 точкам равно ожидаемой площади пика. Разница в последнем (шестом) знаке вызвана округлением отклика до целых чисел. График отклонений от среднего очень напоминает синусоиду. Отклонение измеренной площади от истинной никогда не будет больше, чем амплитуда этой синусоиды. Поэтому естественный следующий шаг – построить зависимость амплитуды синусоиды от частоты измерений.

## Гауссиана: максимальная погрешность в зависимости от частоты оцифровки



Для измерения площади Гауссианы с точностью до 0.1% требуется всего 0.63 точек на  $\sigma$  или 5 отличных от нуля точек на весь пик!

На этом слайде показана зависимость максимальной погрешности интегрирования (абсцисса) от частоты оцифровки, выраженной в измерениях на сигму. По мере роста числа точек максимальная погрешность падает очень быстро, график построен в логарифмическом масштабе. Наиболее интересные уровни погрешности соответствуют одному проценту площади и одной десятой процента. Результаты приятные, по сравнению с оценкой Дайсона позволяют поднять скорость анализа в 8 раз.



На рисунках этого слайда приведена зависимость максимальной погрешности от частоты оцифровки для пиков разной формы. Форма модельных пиков приведена на каждом рисунке; граничные частоты сведены в таблицу.

В реальном анализе вклад ошибок, происходящих от частоты оцифровки, суммируется с ошибками, причиной которых является шум. В большинстве случаев уровень погрешности измерения определяется именно шумом. Приведенную на этом слайде таблицу можно использовать для принятия решения о том, можно ли пренебречь погрешностью, связанной с частотой оцифровки по сравнению с погрешностью, порожденной электронным или химическим шумом. Выведенные закономерности указывают нижний предел частоты измерения и дают определение термина "очень узкий пик". Если принять форму сильно асимметричного пика ЭМГ-3 за худший из приемлемых вариантов пика в хроматографии, то очень узкий пик – это такой, у которого на сигму приходится менее полутора точек.

## Обсуждение

- Функции, по которым вычисляются моменты пика, «пикоподобны», моменты являются интегральными оценками. Площадь – момент нулевого порядка.
- Высоту и другие параметры можно восстановить по известной форме пика
- Восстановление пика при условии известной формы: от 3 точек для Гауссианы, от 4 точек для ЭМГ, измерение площади без ошибок.

Кроме площади, есть другие свойства пика, вычисляемые интегрально. К примеру, дисперсия (квадрат сигмы) – второй центральный момент, и ее можно оценивать по небольшому числу точек аналогично площади.

Есть много других параметров пика, которые нужно оценивать. Интегрирование не может оценить, к примеру, высоту пика. Для определения высоты чрезвычайно полезно иметь представление о форме пика. Если такие представления есть, форма может быть восстановлена аппроксимацией. Для восстановления Гауссова пика потребуется три точки, ЭМГ-четыре.

## ВЫВОДЫ

- ▶ Метод прямоугольников (или трапеций) является лучшим методом измерения площади пиков, обеспечивая ее несмещенную оценку
- ▶ Площадь пика может быть надежно (0.1%) измерена для Гауссовых пиков при 0.63 точек на  $\sigma$ , для несимметричных пиков эта величина может вырасти до 1.6 точек на  $\sigma$

Основной вывод этой презентации - то, что метод прямоугольников является несмещенной оценкой площади пика. Как ни странно, несмотря на его тривиальность, мы не смогли найти его в литературе. Искать можно начиная с конца 19 века...

***Спасибо за внимание!***