

Функции Ферми-Дирака

Н.Н. Калиткин, С.А. Колганов

1 декабря 2019 г.

Оглавление

1	ВВЕДЕНИЕ	5
2	СВОЙСТВА ФУНКЦИЙ ФЕРМИ-ДИРАКА	11
3	ФУНКЦИИ ЦЕЛОГО ИНДЕКСА	27
4	ЭКСПОНЕНЦИАЛЬНО СХОДЯЩИЕСЯ КВАДРАТУ- РЫ	31
5	ФУНКЦИИ ПОЛУЦЕЛОГО ИНДЕКСА	45
6	АППРОКСИМАЦИИ ФУНКЦИЙ ФД	47

Глава 1

ВВЕДЕНИЕ

§1.1. Определение функций

Функции Ферми-Дирака (далее ФД) возникают в задачах квантовой механики при описании свойств вещества, обусловленных поведением электронов или других фермионов. Как известно, функция распределения электронов в 6-мерном пространстве импульсов \mathbf{p} и координат \mathbf{r} имеет следующий вид:

$$f(\mathbf{p}, \mathbf{r}) = \frac{const}{1 + \exp\left(\frac{\frac{1}{2}\mathbf{p}^2 + \phi(\mathbf{r}) - \mu}{T}\right)}. \quad (1.1)$$

Здесь $\phi(\mathbf{r})$ есть потенциал электрического поля, μ – химический потенциал, а T – температура. Константа есть статистический вес частицы, определяемый её спином; для электрона $const = 2$. Все физические формулы написаны в атомной системе единиц; за единицу берутся масса электрона, заряд электрона и постоянная Планка.

При решении квантово-механических задач используются различные моменты фермиевского распределения. Они равны сверткам различных степеней импульса \mathbf{p} с этим распределением по объёму в импульсном пространстве $d\mathbf{p} = 4\pi^2 dp$. Например, электронная плотность в пространстве координат есть свертка нулевой степени импульса:

$$\rho(\mathbf{r}) = \int f(\mathbf{p}) d\mathbf{p} = \int_0^\infty \frac{8\pi^2 dp}{1 + \exp\left[\frac{\frac{1}{2}\mathbf{p}^2 + \phi(\mathbf{r}) - \mu}{T}\right]}. \quad (1.2)$$

Плотность кинетической энергии есть свертка квадрата импульса:

$$E(\mathbf{r}) = \int \left(\frac{1}{2}\mathbf{p}^2 * f(\mathbf{p}) d\mathbf{p}\right) = \int_0^\infty \frac{4\pi p^4 dp}{1 + \exp\left[\left(\frac{1}{2}\mathbf{p}^2 + \phi(\mathbf{r}) - \mu\right)/T\right]}. \quad (1.3)$$

Нахождение проводимости требует нахождения потока частиц, т.е. умножения $f(\mathbf{p}, \mathbf{r})$ на импульс \mathbf{p} ; это эквивалентно введению в числитель

подынтегрального выражения (1.2) степени p^3 вместо p^2 . Теплопроводность определяется через поток кинетической энергии; это эквивалентно введению в числитель подынтегрального выражения (1.3) степени p^5 вместо p^4 .

В таких свертках принято делать следующую замену переменных:

$$t = \frac{p^2}{2T}, x = \frac{\mu - \phi(\mathbf{r})}{T}. \quad (1.4)$$

Тогда различные моменты с точностью до константы приобретают следующий вид:

$$I_k(x) = \int_0^\infty \frac{t^k dt}{1 + \exp(t - x)}, x \in (-\infty; +\infty). \quad (1.5)$$

Например, $k = 1/2$ для электронной плотности (1.2), $k = 1$ для потока частиц, $k = 3/2$ для плотности кинетической энергии (1.3) и $k = 2$ для плотности потока энергии (электронной теплопроводности). Целые k соответствуют нечетным моментам импульса, а полуцелые – четным.

В математической теории функций ФД рассматриваются произвольные индексы k и комплексные значения x .

Интеграл в (1.5) сходится при $k > -1$. Способ доопределения функции ФД для индекса $k < -1$ будет показан далее. Заметим, что при целых отрицательных k функции ФД имеют полюс.

В физических задачах нужны только целые и полуцелые индексы и вещественные значения x . Вычислению функций таких индексов в основном посвящена данная книга. Однако многие приведенные далее выражения справедливы для произвольных индексов k .

В квантово-механических моделях атома возникает ещё одна специфическая функция, связанная с вычислением обменной энергии в квазиклассическом приближении. Пояним подробнее. При квантово-механических расчетах атома широко используется приближение Хартри–Фока, при котором в многоэлектронном уравнении Шредингера многоэлектронная волновая функция заменяется не произведением электронных функций (как в приближении Хартри), а детерминантом, составленным из одноэлектронных функций. Это эффективно учитывает обменное взаимодействие. При этом к электростатическому взаимодействию электронов добавляется дополнительное слагаемое, называемое потенциалом обменного взаимодействия. Однако, такое приближение ещё очень сложно для численных расчетов.

Вычисления обменного потенциала в квазиклассическом приближении дает гораздо более простое выражение. Такое приближение называют приближением Слэтера, который первым написал его для нулевой

температуры. При ненулевой температуре возникает следующая функция:

$$J(x) = \int_{-\infty}^x [I_{-\frac{1}{2}}(\xi)]^2 d\xi. \quad (1.6)$$

Её называют интегральной функцией Ферми-Дирака. Её свойства и способы вычисления также будут рассмотрены в этой книге.

§1.2. История

Свойствам функций ФД посвящено лишь несколько работ. Дадим их обзор.

Функции ФД впервые появились на заре развития квантовой механики в работах Паули [1] и Зоммерфельда [2] при описании частично вырожденного электронного газа в металлах. Основные свойства этих функций изложены в статье МакДугала и Стоунера [3]. В ней приведены сходящийся ряд при $x < 0$ и асимптотическое разложение при $x \rightarrow +\infty$. Заметим, что в асимптотическом разложении в формуле (5.3) имеется опечатка. Вместо сомножителя $(k - r + 2)$ следует писать $(k - 2r + 2)$.

В работе Калиткина [4] построен ряд, сходящийся при любых значениях x . Однако при $x > 3$ скорость сходимости невелика и фактически этот ряд удобен лишь при $-\infty < x < 1$.

Улучшенное выражение для асимптотического разложения получено в работе Glasser [5]; в нем содержится связь между значениями $I_k(x)$ и $I_k(-x)$.

Интегральную функцию ФД ввел Киржниц [6]. Ее разложение в сходящийся ряд при $x < 0$ и асимптотический ряд при $x \rightarrow +\infty$ было построено в работах Каликина, Кузьминой, Луцкого и Колганова.

Ряды, построенные в указанных выше работах, позволяют вычислять функции ФД с высокой точностью при $x < 0$ и при $x > 50$. Остается решить вопрос о практическом вычислении функций ФД в диапазоне $0 < x < 50$.

Функции ФД в указанном диапазоне (и вообще при любом аргументе x) и при заданном k можно найти с точностью ошибок компьютерного округления, непосредственно вычисляя интеграл (1) по каким-либо квадратурным формулам на достаточно подробной сетке. Правда такое вычисление по обычным классическим квадратурным формулам чрезмерно трудоемко и его нецелесообразно использовать как компьютерную подпрограмму. С его помощью можно составить подробные многозначные таблицы этих функций. Первый пример такой таблицы был приведен уже в [3]. Однако такие таблицы имеют огромный объем и также непригодны для создания компьютерных подпрограмм.

Правда, в последние годы для функций ФД полуцелого индекса Ка-

литкиным и Колгановым разработаны экспоненциально сходящиеся квадратуры, радикально уменьшающие трудоемкость вычислений [7-9]. Они уже пригодны для создания подпрограмм, хотя их трудоемкость ощутимо больше, чем для вычислений с помощью упомянутых выше рядов.

Поэтому остается актуальной задача построения быстрых компьютерных подпрограмм. Вот почему ряд работ был посвящен вычислению таблиц функций для конкретных индексов и построению для них несложных экономичных аппроксимаций. В них выбирают некоторый разумный вид аппроксимирующих формул с достаточным числом свободных параметров. Эти параметры подбирают для получения наилучшей точности аппроксимации в некоторой норме. Наиболее разумным представляется получение наилучшей относительной точности в норме S .

Такие аппроксимации для функций целых и полуцелых индексов строились в цитированных работах. Для функций целых индексов аппроксимации с 16 верными десятичными знаками построены в [10], что в полной мере исчерпывает проблему при расчетах с 64-битовыми числами. Для функций полуцелых индексов опубликованные аппроксимации имеют меньшую точность.

В работе МакДугала и Стоунера [3] рассмотрены способы вычисления функций ФД индексов $k = 3/2, 1/2, -1/2$. Но эти способы ещё слишком сложны и трудоемки. Они непригодны для практического использования.

В работе Коди и Тетчера [11] для функций ФД индексов $k = 3/2, 1/2, -1/2$ построен набор аппроксимаций, составленных из трех кусков. Вид этих формул выбран удачно. При $x \leq +1$ аппроксимация является отношением двух многочленов от e^x ; она точно передает главный член асимптотики при $x \rightarrow -\infty$. При $x \geq +4$ аппроксимация есть произведение множителя x^{k+1} на отношение многочленов от x^{-2} . Она точно передает главный член асимптотики при $x \rightarrow +\infty$. В промежутке $+1 \leq x \leq +4$ используется отношение двух многочленов одинаковой степени от x . Наборы с небольшим числом коэффициентов обеспечивают невысокую точность. Для набора с наибольшим числом коэффициентов заявлена относительная точность 10^{-12} . К сожалению, фактическая точность существенно хуже. В Табл.1 приведены относительные рассогласованности смежных формул на стыках. Они на много порядков превышают заявленную погрешность, особенно для $k = 3/2$. По-видимому, в табличных коэффициентах имеются опечатки. Поэтому на практике использовать эти формулы невозможно.

В работе Theiler [12] построены прецизионные аппроксимации для функций ФД индексов $k = 3/2, 1/2, -1/2$. В ней при $x \leq -1$ используется классический ряд по степеням e^x . При $x \geq 30$ использован асимптотический ряд. Промежуток разбивается на 4 отрезка и на каждом отрезке табулированная функция аппроксимируется полиномами Чебышева.

Окончательная относительная погрешность оценивается в 10^{-14} . Сама по себе такая точность хороша. Однако вид аппроксимации выбран менее удачно, чем в работе Коди и Тэтчера. Общее число коэффициентов аппроксимации из-за этого велико. Использование чебышевских многочленов не обеспечивает точность производных. Кроме того, на практике требуются функции и других полуцелых индексов.

Построить единую аппроксимирующую формулу высокой точности вряд ли представляется возможным.

Другой принцип построения формул был предложен в работах Калиткина и Кузьминой [13-15]. В них различные функции ФД выражались через функцию $I_{-1/2}(x)$. Это позволило построить аппроксимацию из 2 кусков. При построении левой части использовались качественные соображения поведения при $x \rightarrow \infty$, а для правой части - асимптотики при $x \rightarrow +\infty$. Это позволило при небольшом числе коэффициентов получить относительную погрешность до $3\Delta 10^{-8}$ для функций полуцелых индексов и $7\Delta 10^{-6}$ для интегральной функции ФД.

§1.3. Погрешность округления

В данной работе рассматриваются методы прямого вычисления функций ФД целых и полуцелых индексов, а также интегральной функции ФД с заданной точностью. Желательно вычислять функции ФД с предельно высокой точностью ϵ , допускаемой компьютером. Поэтому опишем точность, которую дают распространенные сейчас процессоры при вычислении с плавающей точкой.

Все операции с плавающей точкой выполняются арифметическим сопроцессором (*Floating Point Unit*). Его архитектура фактически не развивается с момента принятия в 1985 году стандарта IEEE-754 [ссылка], поэтому предельная разрядность чисел остается на уровне 80 бит. Однако в наиболее распространенных математических обеспечениях для записи в память используются либо 32-битовые числа (*single precision*), либо 64-битовые числа (*double precision*). В обоих случаях не полностью используются возможности линейки процессора. В начале 2000-х годов для языка C++ была возможность использовать 80-битовые числа (*long double precision*); в настоящее время эта возможность не поддерживается и мало где сохранилась. Для суперкомпьютерных вычислений используют 128-битовые числа; но вычисления делаются программными средствами на тех же процессорах. Возможно и вычисление с произвольной разрядностью, но они также выполняются программными средствами. Например, для языка C++ это библиотека *boost::multiprecision*.

При вычислениях с плавающей точкой побитовая запись числа выглядит следующим образом: знак числа – 1 бит, двоичный порядок числа – p бит, мантисса – m бит. Порядок числа может быть положительным или отрицательным; но при записи в память к нему автоматически добавляется положительная константа, равная по модулю максимально

Таблица 1.1: Таблица 1. Компьютерные числа с плавающей точкой

точность	бит	p	m	P	$\lg \varepsilon$
float	32	8	23	± 38	-7.5
double	64	11	52	± 308	-16.2
long double	80	15	64	± 4932	-19.3
Quadruple double	128	15	112	± 4932	-34.3

возможному отрицательному порядку, так что в память записывается неотрицательное число. Это эквивалентно тому, что из p разрядов порядка тратится 1 разряд на знак порядка и $p - 1$ на модуль порядка.

При записи мантииссы также имеется небольшая "хитрость". Первый разряд мантииссы всегда равен 1, поэтому при записи он отбрасывается. При считывании числа в процессор этот разряд автоматически добавляется. Такой прием позволяет фактически удлинять записываемую мантииссу на 1 бит.

Нетрудно посчитать, что максимально возможный порядок в двоичной системе есть $2^{p-1} - 1$; для перевода в десятичную систему его нужно умножить на $\lg 2$. Разрядность процессора длиннее, чем считанное из памяти число, за исключением *long double*. После вычисления результат, записанный на процессоре, оказывается длиннее отведенного в памяти места. Поэтому при записи результата в память производится округление. Ошибка такого округления не превышает половины последнего отброшенного разряда. С учетом "спрятанного" первого разряда мантииссы это означает, что относительная ошибка округления есть $\varepsilon = 2^{-m-2}$.

В случае *long double* длина числа равна линейке процессора и нельзя ни спрятать первый разряд мантииссы, ни иметь лишние разряды для округления; в этом случае $\varepsilon = 2^{-m}$. В Таблице 1 приведены предельные порядки $P = 2^{p-1} \lg 2$ и относительные ошибки единичного округления в форме $\lg \varepsilon$ для рассмотренных выше случаев вычисления. Видно, что для случая *long double* представление числа выбрано не вполне удачно; лучше было бы взять $p = 14, m = 65$.

Наиболее частыми являются 64-битовые вычисления, поэтому мы будем ориентироваться на точность 16^{10} . Вычисление *long double* кажутся заманчивыми, поскольку они повышают точность по сравнению с *double* без увеличения машинного времени. Однако разбиение числа на мантииссу и порядок было недостаточно продумано (лучше было бы перекинуть 1 бит на мантииссу), так что увеличение точности не столь значительно. Но 32-битовыми числами не стоит пренебрегать: именно такую точность дают видеокарты, позволяющие сильно ускорять вычисления за счет конвейерной реализации.

Глава 2

СВОЙСТВА ФУНКЦИЙ ФЕРМИ-ДИРАКА

§2.1. Точное решение

Существует единственный индекс $k = 0$, когда интеграл (1.5) берется в элементарных функциях. Делая в интеграле (1.5) замену переменных $\tau = \exp(t)$, легко получаем

$$I_0(x) = \ln(1 + e^x). \quad (2.1)$$

Для этой функции выполняется важное соотношение, связывающее функцию положительного и отрицательного аргументов:

$$I_0(x) = I_0(-x) + x, x > 0. \quad (2.2)$$

В его справедливости нетрудно убедиться, подставляя (2.1) в (2.2).

Из (2.1) нетрудно получить асимптотическое поведение этой функции: $I_0(x) \approx e^x$ при $x \rightarrow -\infty$, $I_0(x) \approx x$ при $x \rightarrow +\infty$. Видно, что левая и правая асимптоты качественно отличаются друг от друга.

Все прочие функции ФД при $k \neq 0$ в элементарных функциях не выражаются. Для них необходимо разрабатывать алгоритмы, имеющие хорошую точность и экономичность. Эти алгоритмы основываются на различных свойствах функций ФД. Поэтому рассмотрим их основные свойства.

§2.2 Связь функций соседних индексов

Пусть $k > 0$. Вычислим производную функции ФД, дифференцируя (1.5) по :

$$I'_k(x) = \int_0^\infty \frac{d}{dx} \left(\frac{1}{1 + \exp(t - x)} \right) t^k dt \quad (2.3)$$

Поскольку дробь в скобках зависит только от выражения $t-x$, то дифференцирование этой дроби по x эквивалентно дифференцированию этой дроби по t со знаком "минус":

$$I'_k(x) = - \int_0^{\infty} \frac{d}{dt} \left(\frac{1}{1 + \exp(t-x)} \right) t^k dt \quad (2.4)$$

Берем получившийся интеграл по частям:

$$I'_k(x) = - \frac{t^k}{1 + \exp(t-x)} \Big|_0^{\infty} + \int_0^{\infty} \frac{d}{dt} (t^k) \frac{dt}{1 + \exp(t-x)} \quad (2.5)$$

Поскольку $k > 0$, выражение перед интегралом обращается в нуль при $t = 0$ и при $t \rightarrow +\infty$. Раскрывая в последнем интеграле d/dt , получим:

$$I'_k(x) = k \int_0^{\infty} \frac{t^{k-1} dt}{1 + \exp(t-x)} dt = k I_{k-1}(x) \quad (2.6)$$

Это соотношение связывает производную функции ФД с функцией ФД на единицу меньшего индекса.

Соотношение 2.6 было получено для $k > 0$; при этом в правой части стоит функция с индексом больше -1 . Будем считать это соотношение справедливым для любых значений k . Тогда оно доопределяет функции ФД нецелых индексов $k < -1$: оно выражает эти функции через производные функции ФД на единицу большего индекса. Очевидно, для получения функции ФД с индексом $-1 < k < 0$ надо продифференцировать функцию с положительным индексом $0 < k < 1$. Для нахождения функции ФД с индексом $-2 < k < -1$ надо продифференцировать функцию $-1 < k < 0$, т.е. дважды продифференцировать функцию с положительным индексом $0 < k < 1$. Это рекуррентный процесс, т.е. для каждого уменьшения индекса на единицу требуется лишний раз продифференцировать функцию с положительным индексом. Поэтому следует ожидать, что численное нахождение функции ФД с большими отрицательными индексами будет связано с существенной потерей точности.

При $k = 0$ соотношение (2.6) теряет смысл: слева стоит производная от (2), которая положительна и конечна при любом конечном x . Это означает, что $I_{-1}(x) = \infty$ при любом значении x . Соответственно, не существует $I_k(x)$ при любом целом $k < 0$.

§2.3 Ряд для $x < 0$

Пусть $x > 0$. Поскольку в (5) под интегралом $t > 0$, то $\exp x - t < 1$. Тогда в подынтегральном выражении (1) можно преобразовать дробь в

сходящийся ряд:

$$\frac{1}{1+e^{t-x}} = \frac{e^{x-t}}{1+e^{x-t}} = \sum_{n=1}^{\infty} (-1)^{n-1} e^{n(x-t)}. \quad (2.7)$$

Умножим последнее выражение на t^k и проинтегрируем. Интеграл от $t^k \exp(-nt)$ выражается через $(k+1)$. В итоге получаем следующий ряд [Stoner, McAougal]:

$$I_k(x) = \Gamma(k+1) \sum_{n=1}^{\infty} (-1)^{n-1} \frac{e^{nx}}{n^{k+1}}, \quad x < 0. \quad (2.8)$$

Этот ряд сходится при любых $x < 0$, но сходимость неравномерная: она быстро ухудшается при $x \rightarrow -0$. Если $x = 0$, то ряд (5) сходится для индексов $k > -1$, но сходимость при этом крайне медленная. При $x > 0$ ряд расходится для любого k .

Замечание. Почленное дифференцирование ряда (5) с индексом k точно дает ряд для индекса $k-1$.

§2.4. Всюду сходящийся ряд

Этот ряд был предложен в [Калиткин 68]. Применим к подынтегральному выражению (1) следующее преобразование:

$$\frac{1}{1+e^{t-x}} = \frac{2e^{-t}}{(1+2e^{-x}) - (1+2e^{-t})} = \frac{2e^{-t}}{1+2e^{-x}} \left(1 - \frac{1-2e^{-t}}{1+2e^{-x}} \right)^{-1} = \frac{2e^{-t}}{1+2e^{-x}} \sum_{m=0}^{\infty} \left(\frac{1-2e^{-t}}{1+2e^{-x}} \right)^m. \quad (2.9)$$

При допустимых значениях переменных ($-\infty < x < +\infty, 0 < t < +\infty$) разложение в ряд здесь всегда сходится, так как $|(1-2e^{-t})/(1+2e^{-x})| < 1$. Подставляя последнее разложение в интеграл (1), получаем ряд, сходящийся при любых x :

$$I_k(x) = 2\Gamma(k+1) \sum_{n=0}^{\infty} \frac{b_n^{(k)}}{(1+2e^{-x})^{n+1}}, \quad (2.10)$$

где

$$b_n^{(k)} = \frac{1}{\Gamma(k+1)} \int_0^{\infty} (1-2e^{-t})^n e^{-t} t^k dt = \sum_{p=0}^n \frac{(-1)^p 2^p n!}{p!(n-p)!(p+1)^{k+1}}; \quad b_0^{(k)} = 1. \quad (2.11)$$

Сумма в (2.11) получается раскрытием скобки в подынтегральном выражении по формуле бинома и выражением получившихся интегралов через Γ -функцию.

В подынтегральном выражении для $b_n k$ стоит скобка $(1 - 2e^{-t})$, которая по модулю не превосходит 1; поэтому всегда $|b_n^{(k)}| \leq 1$. Следовательно, ряд (8) сходится при любом x не медленней, чем геометрическая прогрессия со знаменателем $(1 + 2e^{-x})^{-1}$. На любом ограниченном справа полубесконечном интервале эта сходимость будет равномерной. Практически при $x \leq 0$, и даже $x \leq 1$, сходимость достаточно быстрая, и ряд (8) удобен для прямых вычислений функций ФД. Этим ряд (8) удобнее ряда (5), который практически пригоден лишь при $x \leq -1$. Однако при $x > 1$ скорость сходимости ряда (8) быстро ухудшается.

Отметим одно качественное отличие ряда (8) от ряда (5). Почленное дифференцирование ряда (8) для индекса k не дает такого же ряда для индекса $k = -1$.

Коэффициенты ряда. Вычисление $b_n^{(k)}$ с помощью суммы (9) легко выполняется для первых членов ряда:

$$b_0^{(k)} = 1; b_1^{(k)} = 1 - 2^{-k}. \quad (2.12)$$

Однако для больших n так вычислять коэффициенты из суммы (9) неудобно. Слагаемые суммы знакопеременны. Члены суммы содержат в себе биномиальные коэффициенты. При больших n центральные коэффициенты на много порядков больше крайних коэффициентов. Поэтому суммирование знакопеременных членов в (9) приводит к потере точности. Уже при $n = 10$ теряется много значащих цифр, а при $n = 20$ потеря точности становится огромной.

Существует способ вычисления $b_n^{(k)}$ без потери точности. Справедливы следующие рекуррентные соотношения, выражающие коэффициенты для индекса k через коэффициенты для индекса $k - 1$:

$$b_0^{(k)} = 1; b_n^{(k)} = \frac{1}{n+1} (b_n^{(k-1)} + n b_{n-1}^{(k)}), n > 1. \quad (2.13)$$

Равенство (11) проверяется интегрированием (9) по частям. Предполагая $k > 0$ и умножая все на $\Gamma(k+1)$, проинтегрируем сомножитель e^{-t} :

$$\Gamma(k+1) b_n^{(k)} = -(1 - 2e^{-t})^n e^{-t} t^k \Big|_0^\infty + \int_0^\infty e^{-t} \frac{d}{dt} [(1 - 2e^{-t})^n t^k] dt. \quad (2.14)$$

Первое слагаемое равно 0. Преобразуем подынтегральное выражение:

$$e^{-t} \frac{d}{dt} [(1 - 2e^{-t})^n t^k] = n e^{-t} (1 - 2e^{-t})^{n-1} (2e^{-t}) t^k + k (1 - 2e^{-t})^n t^{k-1} \quad (2.15)$$

Выражение $(2e^t)$ в первом слагаемом преобразуем так: $(2e^t) = 1 - (1 - 2e^{-t})$. Тогда последний интеграл выражается через коэффициенты следующим образом:

$$n \Gamma(k+1) b_{n-1}^{(k)} - n \Gamma(k+1) b_n^{(k)} + k \Gamma(k) b_n^{(k-1)}. \quad (2.16)$$

Приравнивая полученное выражение величине $\Gamma(k+1)b_n^{(k)}$, доказываем формулу (11).

Поскольку $|b_n^{(k)}| \leq 1$, а коэффициенты в правой части (11) суммируются с положительными весами $1/(n+1)$ и $n/(n+1)$, причем сумма весов равна 1, то вычисления по формуле (11) будут устойчивыми. Заметим, что если для некоторого $k-1$ коэффициенты $|b_n^{(k-1)}| \geq 0$ при всех n , то коэффициенты $b_n^{(k)}$ также будут неотрицательны. Поскольку при увеличении k на 1 вычисления устойчивы, то можно далее снова увеличить k на 1 и повторить этот процесс неограниченное количество раз, не теряя точности.

Заметим, что, рекуррентно суммируя (11) по n , получаем следующее выражение:

$$b_n^{(k)} = \frac{1}{n+1} \sum_{p=0}^n b_p^{(k-1)} \quad (2.17)$$

Вычисление по формуле (12) также устойчиво.

Практически нас интересуют только целые или полуцелые индексы. Для $k=0$ коэффициенты вычисляются точно:

$$b_n^{(0)} = \frac{1 + (-1)^n}{2(n+1)}, n \geq 0; \quad (2.18)$$

эта формула легко получается из интеграла (9) заменой переменных $\tau = (1 - 2e^{-t})$. Заметим, что все нечетные коэффициенты (13) равны 0, а четные положительны. Для других целых индексов $k > 0$ коэффициенты вычисляются рекуррентным применением (11). Видно, что все они будут положительными. В принципе, коэффициенты для целых k выражаются через обыкновенные дроби, но эти выражения чрезмерно громоздки.

Для полуцелых индексов k наиболее удобно начать с индекса $k = -1/2$. В этом случае в интеграле (9) удобно сделать замену $t = \tau^2$. Получившиеся интегралы вычисляются по квадратурным формулам. По найденным коэффициентам и рекуррентным формулам (11) легко вычисляются значения коэффициентов для более высоких полуцелых индексов.

Индекс $k = -3/2$. Он требуется в практических приложениях. Формула (11) очевидным образом обращается в сторону уменьшения индекса:

$$b_0^{(k)} = 1; b_n^{(k-1)} = (n+1)b_n^{(k)} - nb_{n-1}^{(k)}, n > 1. \quad (2.19)$$

Однако эта формула содержит вычитание с большими весами, и вычисления по ней неустойчивы. Выгоднее воспользоваться дифференциальным соотношением (4), полагая в нем $k = -1/2$, и продифференцировать ряд (8) для этого индекса почленно. Это дает

$$I_{-\frac{3}{2}}(x) = -8\sqrt{\pi}e^{-x} \sum_{n=0}^{\infty} \frac{(n+1)b_n^{(-1/2)}}{(1+2e^{-x})^{n+2}} \quad (2.20)$$

Этот ряд также сходится при любых x , а вычисления по этой формуле устойчивы.

Коэффициенты для всех целых и полуцелых индексов от $k = -1/2$ до $k = 4$ приведены в Таблице 1 в виде десятичных дробей. Число значащих цифр рассчитано на получение относительной точности $\varepsilon = 10^{-16}$. Число коэффициентов достаточно для вычисления функций ФД при $x \leq 1$ с указанной точностью. Для значений $x \leq 0$ достаточно суммировать до $n = 35$.

§2.5. Изменение знака аргумента

В [4] получено фундаментальное соотношение, связывающее функции ФД произвольного индекса $k > -1$ от положительного и отрицательного аргументов. Запишем его в несколько преобразованном виде, предполагая $x > 0$:

$$I_k(x) \approx \cos(\pi k) I_k(-x) + \frac{x^{k+1}}{k+1} \left[1 + \sum_{n=1}^N (2 - 2^{2-2n}) \frac{\zeta(2n)}{x^{2n}} \prod_{p=1}^{2n} (k+2-p) \right]; \quad (2.21)$$

здесь ζ — функция от четного аргумента выражается через числа Бернулли:

$$\zeta(2n) = \frac{2^{2n-1} \pi^{2n} |B_{2n}|}{(2n)!}. \quad (2.22)$$

Заметим, что при сравнении данных формул с [4] надо учитывать одно обстоятельство: определение функции ФД в [4] отличается от (1) делителем $\Gamma(k+1)$.

Представление (2.21) является асимптотическим. Это означает, что при фиксированном x и возрастании N , т.е. увеличении числа членов суммы, точность сначала возрастает до некоторого номера $N(x)$, а при дальнейшем добавлении членов начинает ухудшаться. Оптимальное $N(x)$ монотонно возрастает при $x \rightarrow +\infty$. Таким образом, представление (2.21) позволяет получать высокую точность при положительных $x \gg 1$, но при умеренных $x > 0$ его точность невелика.

Нетрудно убедиться, что почленное дифференцирование разложения (2.21) дает разложение для функции ФД с индексом $k-1$. Однако оптимальное число членов $N(x)$ оказывается разным для функций соседних индексов.

В Таблице 2 приведены значения первых 12 чисел Бернулли и ζ — функции. Такого количества достаточно для обеспечения относительной погрешности 10^{16} при аргументах $x \geq 50$.

Таблица 2

Целые индексы. Сам асимптотический ряд (2.21) был получен ещё в ранних работах [3]. Однако там ещё не была обнаружена связь с функциями от отрицательного аргумента. Для функций произвольного ин-

Таблица 2.1: Коэффициенты $b_n^{(k)}$ ряда №

$n \backslash k$	$-1/2$	$1/2$	1
0	1.0000000000000000	1.0000000000000000	1.0000000000000000
1	-0.41421356237309515	0.29289321881345243	0.5000000000000000
2	0.48097395201231308	0.35558679654640596	0.44444444444444442
3	-0.31443745684377605	0.18808073319886048	0.33333333333333331
4	0.35496973905796525	0.22145853437068147	0.30666666666666664
5	-0.2639146991274261	0.14056299545433018	0.25555555555555554
6	0.29305220065772453	0.16234716762624365	0.23945578231292514
7	-0.23206012422784597	0.11304625614448247	0.20952380952380950
8	0.25480909731988405	0.12879768294174931	0.19858906525573192
9	-0.20961507725316636	0.09495640692225774	0.17873015873015871
10	0.22827198241767724	0.10707600469456860	0.17074642529187981
11	-0.19269402581588935	0.08209516881869711	0.15651755651755650
12	0.20850412514677033	0.09181893469008735	0.15039490424105806
13	-0.17934271566143664	0.07245024537926421	0.13965241108098250
14	0.19305732762114028	0.08049071752872261	0.13478669478669478
15	-0.16845640223916253	0.06493152254322979	0.12636252636252635
16	0.1805638658684417	0.07173342509177166	0.12238964418895212
17	-0.15935707527675091	0.05889506396018707	0.11559021951178811
18	0.17019275877281642	0.06475283737137809	0.11227660685050288
19	-0.15160238383284011	0.05393507631116717	0.10666277650797773
20	0.16140657992677773	0.05905276695952958	0.10385117037040963
21	-0.14488973151960657	0.04978265339229612	0.09913066262630009
22	0.15384051431503362	0.05430690821502384	0.09671099298470294
23	-0.13900414944762135	0.04625228081241362	0.09268136827700697
24	0.1472371586201546	0.05029167592472326	0.09057411354592670
25	-0.13378812964569098	0.04321168340278425	0.08709049379416028
26	0.14140906612540924	0.04684862350362221	0.08523666206241499
27	-0.12912325802251828	0.04056391344911719	0.08219249556018589
28	0.13621609350926644	0.04386226448567407	0.08054733221744463
29	-0.12491855193815647	0.03823623727154638	0.07786242114352980
30	0.13155110392581695	0.04124639426039382	0.07639131286522656
31	-0.12110278860777168	0.03617298229576366	0.07400408433818823
32	0.12733061461119022	0.03893533479017052	0.07267980997348647
33	-0.11761930413804721	0.03433078658639941	0.07054216850367805
34	0.12348848883008756	0.03687814950764765	0.06934300450561377

Таблица 2.2: Коэффициенты $b_n^{(k)}$ ряда №

$n \backslash k$	$-1/2$	$1/2$	1
35	-0.11442237441589655	0.03267535717643808	0.06741680993601339
36	0.11997156715133331	0.03503471420278660	0.06632519418171647
37	-0.11147463877385461	0.03117920491392763	0.06457979433482919
38	0.11673657115056887	0.03337298353537996	0.06358136436832140
39	-0.10874522884795892	0.02982002822579649	0.06199183025911337
40	0.11374786384105723	0.03186704860665650	0.06107471839675545
41	-0.10620838520960593	0.02857953827769788	0.05962055843492794
42	0.11097580026521373	0.03049573041694243	0.05877486670280144
43	-0.10384241822419096	0.02744259067509849	0.05743907427773778
44	0.10839549296611649	0.02924154405934333	0.05665647756539299
45	-0.10162891623985279	0.02639653405283907	0.05542481500962358
46	0.10598587447494355	0.02808992427458597	0.05469825715292267
47	-0.09955213438673234	0.02543071471914185	0.05355871012890345
48	0.10372897595078695	0.02702863841774685	0.05288216835617696
49	-0.09759851721451281	0.02453609530510165	0.05182452498905343
50	0.10160936547917465	0.02604733589675014	0.05119282534490051
51	-0.09575632184163210	0.02370495786331971	0.05020834793442165
52	0.09961370590635248	0.02513719839243354	0.04961701919090997
53	-0.09401531751504925	0.02293067032007274	0.04869818550218941
54	0.09773040328255238	0.02429066546484510	0.04814334246000084
55	-0.09236654391281224	0.02220750101167265	0.04728363991607226
56	0.09594932477105972	0.02350121721797768	0.04676188938507314
57	-0.09080211507512911	0.02153047010947584	0.04595564991291671
58	0.09426157037285680	0.02276320062241450	0.04546401436426342
59	-0.08931505911688105	0.02089522962675958	0.04470628079152570
60	0.09265928674694333	0.02207168957954948	0.04424213590351627
61	-0.08789918624863347	0.02029796577586911	0.04352855306636278
62	0.09113551425473628	0.02142237130727970	0.04308957628551600
63	-0.08654897937862024	0.01973531895281251	0.04241630165605482
64	0.08968406044561429	0.02081145343731716	0.04200042955957113
65	-0.08525950286539473	0.01920431773576092	0.04136405941472915
66	0.08829939474918186	0.02023558754193138	0.04096945215681271
67	-0.08402632596474541	0.01870232410800966	0.04036696021333017

Таблица 2.3: Значения чисел Бернулли и ζ -функции

$2n$	B_{2n}	$\zeta(2n)/\pi^{2n}$
2	$\frac{1}{6}$	$\frac{1}{6}$
4	$-\frac{1}{30}$	$\frac{1}{90}$
6	$\frac{1}{42}$	$\frac{1}{945}$
8	$-\frac{1}{30}$	$\frac{1}{9450}$
10	$-\frac{5}{66}$	$\frac{1}{93555}$
12	$-\frac{691}{2730}$	$\frac{691}{638512875}$
14	$\frac{7}{6}$	$\frac{2}{18243225}$
16	$-\frac{3617}{510}$	$\frac{3617}{325641566250}$
18	$\frac{43867}{798}$	$\frac{43867}{38979295480125}$
20	$-\frac{174611}{330}$	$\frac{174611}{1531329465290625}$
22	$\frac{854513}{138}$	$\frac{155366}{13447856940643124}$
24	$-\frac{236364091}{2730}$	$\frac{236364091}{201919571963756511232}$

декса последнее обстоятельство несущественно, поскольку точность асимптотического ряда хороша лишь при больших значениях x , когда величина $I_k(-x)$ весьма мала.

Однако для целых индексов $k \geq 0$ ситуация кардинально меняется. В этом случае в произведении по p сомножитель со значением $p = 2k + 2$ равен нулю. Тем самым, сумма по p содержит лишь конечное число слагаемых вплоть до $N = 1 + [k/2]$, где квадратные скобки обозначают целую часть числа. Поэтому соотношение (2.21) принимает следующий вид:

$$I_k(x) = (-1)^k I_k(-x) + \frac{x^{k+1}}{k+1} \left[1 + \sum_{n=1}^{1+[k/2]} (2 - 2^{2-2n}) \frac{\zeta(2n)}{x^{2n}} \prod_{p=1}^{2n} (k+2-p) \right]; \quad (2.23)$$

Самое поразительное то, что для целых индексов соотношение (2.23) является не асимптотическим, а точным! В этом можно убедиться следующим образом. Для функции $I_0(x)$ оно совпадает с соотношением (3), которое является точным. Полагая $k = 1$ в (4), подставляя в правую часть (3) и интегрируя, получим для функции $I_1(x)$ соотношение между функциями положительного и отрицательного аргументов, которое также оказывается точным. Последовательно повышая таким интегрированием индекс k на единицу, убедимся в том, что соотношение (2.23) является точным для всех функций целого индекса.

Приведем такие соотношения для нескольких первых индексов:

$$\begin{aligned} I_0(x) &= I_0(-x) + x, x \geq 0; \\ I_1(x) &= -I_1(-x) + \frac{x^2}{2} + \frac{\pi^2}{6}, x \geq 0; \\ I_2(x) &= I_2(-x) + \frac{x^3}{3} + \frac{\pi^2}{3}x, x \geq 0; \\ I_3(x) &= -I_3(-x) + \frac{x^4}{4} + \frac{\pi^2}{2}x^2 + \frac{7\pi^4}{60}, x \geq 0; \\ I_4(x) &= I_4(-x) + \frac{x^5}{5} + \frac{2\pi^2}{3}x^3 + \frac{7\pi^4}{15}x, x \geq 0. \end{aligned} \quad (2.24)$$

Этих значений k достаточно для практических приложений.

§2.6. Интегральная функция ФД

Эта функция определена соотношением (1.6). Исследуем свойства этой функции.

Разложение при $x < 0$. Запишем разложение (5) для индекса $k = -1/2$:

$$I_{-1/2}(x) = \sqrt{\pi} \sum_{n=1}^{\infty} (-1)^{n-1} \frac{e^{nx}}{\sqrt{n}}, x \leq 0. \quad (2.25)$$

Возводя этот ряд в квадрат и группируя члены с одинаковыми экспонентами, получим

$$[I_{-\frac{1}{2}}(x)]^2 = \pi \sum_{n=2}^{\infty} (-1)^n e^{nx} \sum_{p=1}^{n-1} \frac{1}{\sqrt{p(n-p)}}, x \leq 0. \quad (2.26)$$

Почленно интегрируя полученный ряд от $-\infty$ до x , получим искомое разложение:

$$J(x) = \pi \sum_{n=2}^{\infty} (-1)^n a_n e^{nx}, a_n = \frac{1}{n} \sum_{p=1}^{n-1} \frac{1}{\sqrt{p(n-p)}}, x \leq 0. \quad (2.27)$$

Это разложение сходится при $x < 0$ и расходится при $x > 0$. Сходимость тем быстрее, чем больше модуль x . Ряд (2.27) знакопеременный. Фактически он пригоден для вычисления функции (1.6) при $x \leq -1$, как и ряд (5). Значения коэффициентов a_n приведены в Таблице 3.

Всюду сходящийся ряд. Напишем всюду сходящийся ряд (8) для $k = -1/2$:

$$I_{-\frac{1}{2}} = 2\sqrt{\pi} \sum_{n=0}^{\infty} \frac{b_n^{(-1/2)}}{(1+2e^{-x})^{n+1}}, -\infty \leq x \leq +\infty. \quad (2.28)$$

Возводя этот ряд в квадрат и группируя члены с одинаковыми знаменателями, получим

$$[I_{-\frac{1}{2}}]^2 = 4\pi \sum_{n=0}^{\infty} \frac{1}{(1+2e^{-x})^{n+2}} \sum_{p=0}^n b_p^{(-1/2)} b_{n-p}^{(-1/2)}. \quad (2.29)$$

Для получения $J(x)$ надо проинтегрировать (2.29) от $-\infty$ до x . Будем искать интеграл в виде аналогичного разложения с неизвестными коэффициентами c_n :

$$J(x) = 4\pi \sum_{n=0}^{\infty} \frac{c_n}{(1+2e^{-x})^{n+2}}, -\infty \leq x \leq +\infty. \quad (2.30)$$

Продифференцируем (2.30) по x и проведем следующие преобразования:

$$\begin{aligned} J'(x) &= 4\pi \sum_{n=0}^{\infty} \frac{2e^{-x}(n+2)c_n}{(1+2e^{-x})^{n+3}} = 4\pi \sum_{n=0}^{\infty} \frac{[(1+2e^{-x})-1](n+2)c_n}{(1+2e^{-x})^{n+3}} = \\ &= 4\pi \left[\sum_{n=0}^{\infty} \frac{(n+2)c_n}{(1+2e^{-x})^{n+2}} - \sum_{n=0}^{\infty} \frac{(n+2)c_n}{(1+2e^{-x})^{n+3}} \right]. \end{aligned} \quad (2.31)$$

Сдвигая во второй сумме индекс n на 1, получим

$$J'(x) = 4\pi \left[\sum_{n=0}^{\infty} \frac{(n+2)c_n}{(1+2e^{-x})^{n+2}} - \sum_{n=1}^{\infty} \frac{(n+1)c_{n-1}}{(1+2e^{-x})^{n+2}} \right] \quad (2.32)$$

Сопоставляя последнее выражение с (2.29) и приравнявая коэффициенты при одинаковых степенях знаменателя, найдем следующее соотношение между коэффициентами:

$$\begin{aligned} n = 0 : 2c_0 &= [b_0^{(-1/2)}]^2 = 1; \\ n > 0 : (n+2)c_n - (n+1)c_{n-1} &= \sum_{p=0}^n b_p^{(-1/2)} b_{n-p}^{(-1/2)}. \end{aligned} \quad (2.33)$$

Таким образом, коэффициенты формулы (2.30) определяются из рекуррентных соотношений

$$c_0 = \frac{1}{2}; c_n = \frac{1}{n+2} \left[(n+1)c_{n-1} + \sum_{p=0}^n b_p^{(-1/2)} b_{n-p}^{(-1/2)} \right], n > 0. \quad (2.34)$$

Значения коэффициентов c_n приведены в Таблице 3.

Формулы (2.30) и (2.34) определяют ряд, сходящийся при любых значениях x . Скорость его сходимости достаточно хороша при $x \leq 0$, удовлетворительна при $x \leq 1$, но быстро ухудшается при $x > 1$.

Разложение при $x \rightarrow +\infty$. Запишем асимптотическое разложение (15) для произвольного k в более удобном виде:

$$\begin{aligned} I_k(x) &\approx \cos(\pi k) I_k(-x) + \frac{x^{k+1}}{k+1} \sum_{n=0}^N \frac{A_n^{(k)}}{x^{2n}} \\ A_0^{(k)} &= 1; A_n^{(k)} = (2 - 2^{2-2n}) \zeta(2n) \prod_{p=1}^{2n} (k+2-p), n \geq 1. \end{aligned} \quad (2.35)$$

Положим здесь $k = -1/2$; тогда $\cos(\pi k) = 0$, и остается только сумма. Возведем эту сумму в квадрат и опять сгруппируем члены по одинаковым степеням. Суммирование по n идет не до бесконечности, а до N . Однако мы будем пользоваться разложением (2.35) только в том случае, если последние члены суммы достаточно малы по сравнению с главными. Тогда можно приближенно записать

$$[I_{-1/2}(x)]^2 \approx 4x \sum_{n=0}^N \frac{C_n}{x^{2n}}, C_n = \sum_{q=0}^n A_q^{(-1/2)} A_{n-q}^{(-1/2)}, x \rightarrow +\infty. \quad (2.36)$$

Таблица 2.4: Коэффициенты разложений функции $J(x)$

n	a_n	c_n	C_n
0	0.5000000000000000	0.5000000000000000	1.0000000000000000
1	0.47140452079103173	0.05719095841793650	-0.82246703342411309
2	0.41367513459481287	0.32627341363306145	-3.38226010534730559
3	0.36329931618554523	0.05555337454026419	-56.7486676763200464
4	0.32247788425329943	0.24658846860286468	-2076.43981697169329
5	0.28947176887871823	0.05073748578641944	-133516.623919083009
6	0.26245962433780035	0.20002927599276679	-13363920.4954685569
7	0.24002748681137623	0.04624864575737019	-1924202279.42978835
8	0.22113507376025146	0.16919747074124367	-376996608458.572022
9	0.20502010799708259	0.04243339943502982	-96469021655492.7344
10	0.19111818527221525	0.14714269473929523	-31243036135798104.0
11	0.17900522098810093	0.03922350305150957	-12492545181655248896.0
12	0.16835758850494667	0.13051578228471081	-6044381261816933646336.0
13	0.15892459422501767	0.03650451326932291	
14	0.15050930279877994	0.11749388727681387	
15	0.14295500139265471	0.03417682128363331	
16	0.13613550394070195	0.10699573032743079	
17	0.12994810323070083	0.03216234704554120	
18	0.12430837457100294	0.09833733348435542	
19	0.11914629284901028	0.03040122425575591	
20	0.11440329429976415	0.09106388018262407	
21	0.11003002691122073	0.02884749082202021	
22	0.10598460916982465	0.08486060243799427	
23	0.10223126852657950	0.02746553440048496	
24	0.09873926667442291	0.07950240515223343	
25	0.09548204372438338	0.02622743720392143	
26	0.09243653108215864	0.07482385935080225	
27	0.08958259552832193	0.02511104569452402	
28	0.08690258621471895	0.07070053682482372	
29	0.08438096303764077	0.02409857188256913	
30	0.08200398984235426	0.06703696772323861	
31	0.07975947964440647	0.02317557364450819	
32		0.06375862474843785	
33		0.02233020393904351	
34		0.06080644154622106	

Почленно проинтегрируем (2.36) по x от $-\infty$ до x , учитывая значения $C_0 = 1, A_1(-1/2) = -\pi^2/24, C_1 = -\pi^2/12$. Получим следующую сумму:

$$J(x) \approx 2x^2 \left[1 - \frac{\pi^2}{6}(\ln x - j) \frac{1}{x^2} - 4 \sum_{n=2}^N \frac{(n-1)C_n}{x^{2n}} \right], x \rightarrow +\infty, \quad (2.37)$$

где j есть константа, возникающая при интегрировании. Значение этой константы приведено в [6-7]:

$$j = \frac{\pi^2}{2} \left(1 - \frac{2}{3} \ln 2 - \frac{C}{3} \right) - \sum_{n=2}^{\infty} \frac{\ln n}{n^2}, \quad (2.38)$$

где $\approx 0.5772156649015325\dots$ – константа Эйлера.

Медленно сходящуюся сумму (2.38) необходимо вычислить с точностью $\varepsilon \sim 10^{-16}$. Непосредственное суммирование на компьютере требует неприемлемо большого числа членов ряда. Поэтому воспользуемся следующим приемом. Разобьем бесконечную сумму на две:

$$\sum_{n=2}^{\infty} \frac{\ln n}{n^2} = \sum_{n=2}^N \frac{\ln n}{n^2} + \sum_{n=N+1}^{\infty} \frac{\ln n}{n^2}, N \gg 1. \quad (2.39)$$

Первую сумму вычислим непосредственно; при этом суммировать будем с последнего члена, т.к. проведение суммирования в порядке увеличения слагаемых уменьшает ошибки округления. Вторую сумму рассмотрим, как квадратуру средних для интеграла от функции $n^{-2} \ln(n)$ в пределах $N + 1/2 \leq n \leq +\infty$ с шагом $\Delta n = 1$. Сам интеграл легко вычисляется точно заменой переменных $\xi = \ln n$ и равен:

$$\int_{N+1/2}^{\infty} \frac{\ln(n)}{n^2} dn = \frac{1 + \ln(N + 1/2)}{N + 1/2}. \quad (2.40)$$

Для повышения точности добавим к формуле средних поправки Эйлера-Маклорена, содержащие первую и третью производные подынтегральной функции на левой границе (очевидно, эти поправки на правой границе обращаются в нуль). Получим следующее выражение:

$$\sum_{n=N+1/2}^{\infty} \frac{\ln(n)}{n^2} \approx \frac{1 + \ln(N + 1/2)}{N + 1/2} - \frac{2 \ln(N + 1/2) - 1}{24(N + 1/2)^3} + \frac{7[24 \ln(N + 1/2) - 26]}{5760(N + 1/2)^5}. \quad (2.41)$$

Следующая поправка Эйлера-Маклорена есть $O(N^{-7})$; чтобы она не превышала 10^{-16} , достаточно взять $N = 300$. Численный расчет с этими значениями даёт

$$\begin{aligned} \sum \ln n &= 0.93754825431584388, \\ j &= 0.76740941382814898. \end{aligned} \quad (2.42)$$

Коэффициенты C_n для $n \geq 2$ приведены в Таблице 3. Разложение (2.37) имеет асимптотическую сходимость, так что суммировать по n можно только до тех пор, пока члены суммы достаточно быстро убывают. Определение оптимального числа членов N является самостоятельной проблемой.

Глава 3

ФУНКЦИИ ЦЕЛОГО ИНДЕКСА

§3.1. Нулевой индекс.

Напомним, что функции ФД $I_k(x)$ целого индекса k существуют при $k \geq 0$. При $k = 0$ функция ФД выражается по формуле (2.1): $I_0(x) = \log(1 + e^x)$. При целых $k < 0$ функция ФД не существует, так как интеграл (1.5) расходится.

§3.2 Отрицательные аргументы

Классический ряд (2) пригоден для практического вычисления функций при $x \leq -1$. Этот ряд знакопеременный, так что ошибка не превышает первого отброшенного слагаемого. Обычно требуется обеспечить некоторую не абсолютную, а относительную точность ε . Главный член суммы в (5) есть e^x . Если в сумме оставляется N слагаемых, то относительная величина отброшенного члена есть $e^{-Nx}/(N+1)^{k+1}$. Поэтому для заданного ε надо выбрать N из следующего условия:

$$\frac{e^{-Nx}}{(N+1)^{k+1}} \leq \varepsilon. \quad (3.1)$$

Можно определить N из уравнения 3.1 каким-либо итерационным процессом. Но при написании программы можно поступить проще: заранее составить таблицу границ x_N удовлетворяющих уравнению 3.1, и далее отслеживать попадание требуемого значения в эти границы.

Когда N определено, то суммировать отрезок ряда удобнее всего по схеме Горнера:

$$I_k(x) = \Gamma(k+1)e^x \left(\frac{1}{1^{k+1}} - e^x \left(\frac{1}{2^{k+1}} - e^x \left(\frac{1}{3^{k+1}} - e^x \left(\dots - e^x \left(\frac{1}{N^{k+1}} \right) \right) \right) \right) \right). \quad (3.2)$$

Для 64-битовых вычислений следует полагать $\varepsilon = 10^{-16}$. Тогда для $x = -1$ число слагаемых составляет $N = 37 - 40$; оно слабо зависит от k и слегка убывает при увеличении k . При $x \rightarrow -\infty$ число слагаемых N быстро убывает.

Таким образом, это достаточно экономичный способ прямого вычисления функций ФД при $x \leq -1$. Однако, при этом остается нерешенной ещё одна проблема: как вычислять функции на отрезке $-1 \leq x \leq 0$?

Всюду сходящийся ряд. Для произвольных, необязательно целых, индексов k существует всюду сходящийся ряд 2.10. Коэффициенты этого ряда вычислены и приведены в Табл 2.1. Этот ряд сходится при любом x не хуже, чем геометрическая прогрессия со знаменателем $g = (1 + 2e^{-x})^{-1}$. Эта сходимостъ неравномерная, и при $x \rightarrow +\infty$ становится очень медленной. Однако при $x \leq 0$ знаменатель $g \leq 1/3$, и сходимостъ становится достаточно быстрой. Таким образом, ряд 2.10 удобен для вычисления функций ФД при $x \leq 0$.

Ограничимся конечным числом членов ряда (3) и запишем полученную сумму по схеме Горнера:

$$//gh \quad (3.3)$$

напомним, что для целого индекса $\Gamma(k+1) = k!$. Значения коэффициентов $b_n^{(k)}$ можно непосредственно рассчитывать по формулам (5-6) или брать из Таблицы 2.1, где они приведены с 16-ю десятичными знаками. Такого числа знаков достаточно для вычисления с относительной точностью $\varepsilon = 10^{-16}$ (double precision). Описанный алгоритм прост и экономичен. Поскольку коэффициенты b_n^k , то нигде не возникает вычитаний, и относительные ошибки округления могут накапливаться лишь незначительно, практически не ухудшая точность.

Определим требуемое число членов. Нас интересует относительная погрешность. Поэтому общий множитель $2\Gamma(k+1)g$ можно откинуть. Главный член оставшегося произведения есть $b_0^{(k)} = 1$ и все остальные члены можно сравнивать непосредственно с ним. Все коэффициенты $b_n^{(k)} < 1$, положительны и убывают с увеличением номера n . Это означает, что отброшенная часть ряда сходится быстрее, чем геометрическая прогрессия с первым членом g^{N+1} и знаменателем g . Величина знаменателя возрастает от 0 до $1/3$ при возрастании от $-\infty$ до 0. Поэтому сумма отброшенных членов не превышает сумму этой геометрической прогрессии $g^{N+1}/(1-g)$. Нужно, чтобы эта величина не превышала ε .

Эту оценку нетрудно немного усилить. На самом деле величину откинутых членов ряда надо сравнивать не с $b_0^{(k)} = 1$, а с суммой учтенных членов ряда. Их также можно оценить, как сумму геометрической прогрессии с первым членом 1 и знаменателем g . Тогда для сравнения вместо ε надо брать величину $\varepsilon/(1-g)$. Это дает следующую оценку:

$$N \geq \frac{\ln \varepsilon}{\ln g} - 1. \quad (3.4)$$

Эта формула дает нецелое значение N . И его надо округлить вверх до целого:

$$N = \left\lceil \frac{\ln \varepsilon}{\ln g} \right\rceil, \quad (3.5)$$

где квадратные скобки означают целую часть числа. Даже эта оценка завышена, т.к. она не учитывает заметного убывания $b_n^{(k)}$ при возрастании номера n . Однако учет этого эффекта заметно усложнил бы алгоритм, что нецелесообразно.

Наиболее медленная сходимость будет при $x = 0$, когда $\varepsilon = 10^{-16}$; для $g = 1/3$ это дает $N = 33$. При $x \rightarrow -\infty$ величина g быстро стремится к 0; при этом число членов N быстро убывает.

§3.3. Положительный аргумент

Для функций ФД существует точное соотношение, связывающее значения функции при положительном и отрицательном аргументах (2.23). Для индексов $k \leq 4$ выше были приведены конкретные реализации (2.24) общей формулы. Поскольку для $x \leq 0$ алгоритм вычисления построен выше, эти формулы полностью решают вопрос вычисления функций с целым индексом.

Однако построенные выше ряды для отрицательных значений аргумента при $x \rightarrow 0$ требуют суммирования большого числа членов ряда, т.е. сравнительно трудоемки. В Главе № будут построены аппроксимации, существенно снижающие трудоемкость при $x \leq 0$.

Обсудим точность вычислений при $x > 0$. Функции отрицательного аргумента мы вычисляем с относительной точностью ε . Полиномиальная часть формул (2.23-2.24) вычисляется практически точно: относительная ошибка не превышает ошибки единичного округления. Сами функции ФД при $k > -1$ являются монотонно возрастающими по x . Поэтому полиномиальная часть формул (2.23-2.24) превышает вклад функций отрицательного аргумента тем сильнее, чем больше x . Тем самым, при $x > 0$ относительная погрешность описанного алгоритма будет заведомо меньше ε .

Глава 4

ЭКСПОНЕНЦИАЛЬНО СХОДЯЩИЕСЯ КВАДРАТУРЫ

§4.1. Проблема трудоемкости квадратур

Прямое вычисление функций ФД нецелого индекса требует применения квадратурных формул к интегралу (1.5). При произвольных индексах k (в том числе и целых!) такое интегрирование для получения высокой точности (16 верных десятичных знаков) требует очень подробной сетки ($10^5 - 10^6$ узлов) и является чрезмерно трудоемким. Однако для полуцелых индексов k возможно кардинальное уменьшение трудоемкости при специальном преобразовании подынтегрального выражения. Можно построить квадратуры с экспоненциальной, то есть очень быстрой сходимостью.

Квадратуры с экспоненциальной сходимостью были предложены и исследованы в [ссылка на работу 2014 года, наши работы]. В [ссылка на статью SIAM] было рассмотрено интегрирование периодической функции $u(x)$ на ее периоде $[0, 2\pi]$ по формуле трапеций на равномерной сетке с N интервалами. Была доказана следующая

Теорема 1. Пусть $u(x)$ есть функция, аналитическая в полуплоскости $\text{Im} x \geq -a, a > 0$. Тогда справедлива следующая мажорантная оценка погрешности формулы трапеций:

$$|\delta| \leq \frac{2\pi M}{[\exp(aN) - 1]}, \quad (4.1)$$

где $N = \max|u(x)|$ на отрезке интегрирования. TODO:добавить жирный кружок в конце Th. Дадим комментарий. Неаналитичность функции $u(x)$ ниже указанной полуплоскости означает, что на линии $\text{Im}(x) = -a$ лежит по меньшей мере одна особая точка функции. Тогда, в силу перио-

дичности $u(x)$, эти особые точки будут расположены на линии $Im x = -a$ также с периодом 2π . Тогда, одна из этих особых точек будет лежать непосредственно под вещественным отрезком интегрирования. Тем самым, a есть расстояние от отрезка интегрирования до ближайшей особой точки.

Интуитивно ясно, что теорему можно усилить. По-видимому достаточно требовать аналитичности $u(x)$ в полосе $-a \leq Im x \leq a$. Кроме того, для $u(x) = const$, погрешность квадратур трапеций равна 0. Это значит, что вычитая из $u(x)$ константу, мы не меняем фактической погрешности. Поэтому оценку постоянной M можно немного улучшить. Например, $u(x)$ вещественно на вещественной оси, то $M = |\max u + \min u|/2$ на вещественном отрезке интегрирования.

Экспоненциальная сходимость гораздо быстрее степенной. Поэтому такие квадратуры могут обеспечить экономичность расчетов. Покажем, как можно построить экспоненциально сходящиеся квадратуры для функций ФД полуцелых индексов.

§4.2. Сходимость квадратур

Рассмотрим задачу вычисления интегралов от функций $u(x)$, имеющих сколь угодно высокие непрерывные производные на отрезке интегрирования $[a, b]$ (см. напр. [3-5]). Чаще всего на практике берут равномерные или сводящиеся к равномерным сетки $\omega_N = x_n, n = 0, \dots, N$ и используют простейшие квадратурные формулы трапеций, средних, Симпсона и т.п. Погрешность подобных формул имеет оценку $const$, где p есть порядок точности формулы, h – шаг интегрирования, *TODO* Такая сходимость называется **степенной**, поскольку погрешность выражается через степень шага. Она достаточно медленная, и для получения высокой точности требуется большое N . Такие квадратуры довольно трудоемки.

Квадратуры Гаусса-Кристоффеля дают гораздо более быструю сходимость. Например, классическая формула Гаусса для интегрирования на отрезке $[-1, 1]$ с весом $\rho(x) = 1$ имеет погрешность (после упрощения факториальных множителей)

$$\delta \leq \sqrt{\frac{\pi}{N}} \frac{b-a}{4} \left(e \frac{b-a}{8N} \right)^{2N} M_{2N}. \quad (4.2)$$

Квадратура Эрмита для отрезка $[-1, 1]$ с весом $\rho(x) = (1-x^2)^{-1/2}$ имеет погрешность

$$\delta \leq \sqrt{\frac{\pi}{N}} \left(\frac{e}{2\sqrt{2}N} \right)^{2N} M_{2N}. \quad (4.3)$$

Погрешности (11-12) с точностью до логарифмически малых членов можно записать в следующем виде:

$$\delta \sim \alpha \exp(-\beta N). \quad (4.4)$$

Зависимость от числа узлов является не степенной, а экспоненциальной, поэтому такую сходимость будем называть **экспоненциальной**.

Заметим, что формулы со степенной сходимостью порядка p требуют существования непрерывной производной лишь p -ого порядка подынтегральной функции независимо от числа узлов сетки N . Формулы Гаусса-Кристоффеля с N узлами требуют существования $2N$ -й непрерывной производной, т.е. при увеличении N надо соответственно повышать гладкость функции.

Трудоемкость формул Гаусса-Кристоффеля несравненно меньше, чем у квадратур со степенной сходимостью. Однако узлы и веса этих квадратур найдены лишь для отдельных отрезков и весов $\rho(x)$ интегрирования. При этом только для квадратур Эрмита эти веса и узлы найдены в виде простых формул для произвольных N . Для остальных случаев узлы и веса точно вычисляются (через радикалы) лишь для $N \leq 3$ или $N = 5$. Это сильно ограничивает возможности практического использования таких квадратур.

Далее покажем, что если $u(x)$ четно продолжается через обе границы отрезка, то формула трапеций на равномерной сетке дает экспоненциальную сходимость. При этом коэффициент β в экспоненте определяется расстоянием до ближайшей особой точки в комплексной плоскости. Это открывает новые возможности для построения квадратур малой трудоемкости.

§4.3. Случай экспоненциальной сходимости

Пусть $u^{(p)}(x)$ существуют и непрерывны на $[a; b]$ при любых p . Требуется вычислить интеграл

$$U = \int_a^b u(x) dx. \quad (4.5)$$

Введем равномерную сетку ω_n с $x_0 = a, x_N = b$ и воспользуемся формулой Эйлера-Маклорена, базирующейся на формуле трапеций [6, 7]:

$$U_N = h\left(\frac{u_0}{2} + u_1 + u_2 + \dots + u_{N-1} + \frac{u_N}{2}\right) + \sum_{p=1} (-1)^p a_p h^{2p} (u_N^{(2p-1)} - u_0^{(2p-1)}), \quad a_p \sim M_{2p-1}. \quad (4.6)$$

Если оборвать эту сумму на члене P , то первый отброшенный член будет остаточным. Его величина есть $\delta_P = O(h^{2P+2})$. В этом случае формула (4.6) имеет степенную сходимость.

Пусть $u(x)$ такова, что все её нечётные производные на правой и левой границах одинаковы: $u^{(2p-1)}(a) = u^{(2p-1)}(b)$. Тогда в (4.6) сумма обращается в нуль. Оставшаяся часть квадратур является просто формулой трапеций. Из этого следуют

Утверждение 1. Пусть подынтегральная функция $u(x)$ имеет сколько угодно высокие производные, причем нечетные производные на правой и левой границах одинаковы: $u^{(2p-1)}(a) = u^{(2p-1)}(b)$. Тогда формула трапеций на равномерной сетке имеет сходимость выше степенной. \square

Частный случай. Утверждение 1 справедливо, если $u(x)$ чётно продолжается через обе границы отрезка: $u^{(2p-1)}(a) = u^{(2p-1)}(b) = 0$. \square

Таким образом, установлен класс функций, для которого формула трапеций имеет сверхстепенную сходимость. Остается найти закон этой сходимости. Проведем ее изучение на следующем тестовом примере:

$$U(q, r, c) = \int_0^\pi \frac{(c^2 - 1)c^x \cos(rx)}{(c^2 - 2c \cos x + 1)^q} dx, c > 1. \quad (4.7)$$

Параметры $r \geq 0$, $q \geq 1$ берутся целыми. Тогда подынтегральное выражение чётно на обеих границах отрезка, его нечетные производные на границах обращаются в нуль, и пример удовлетворяет требованиям Утверждения 1. При $q = 1$ известно точное значение интеграла [8]:

$$U(1, r, c) = \pi. \quad (4.8)$$

При $q \neq 1$ интеграл (4.7) не выражается через элементарные функции от параметров.

Для тщательного численного выявления закономерностей все расчеты проводились с повышенной разрядностью (45 десятичных знаков) с помощью библиотеки языка *C++ boost::multiprecision*.

Расчеты интеграла (4.7) при фиксированных параметрах проводились на сетках с разным числом интервалов N . Погрешность расчетов при $q = 1$ определялась непосредственным сравнением с точным ответом (4.8). На рис.1 показана зависимость погрешности от N при $r = 0$ и различных значениях c в полулогарифмическом масштабе. Каждому значению c соответствует своя линия погрешности. Видно, что при всех значениях c кривые погрешности в этом масштабе являются прямыми. Это означает, что погрешность подчиняется закону

$$\ln \delta_N = \alpha - \beta N, \beta = \text{const} \cdot \ln c. \quad (4.9)$$

При других значениях параметров картина была аналогичной. На рис.2 показан случай $q = 1, r = 2$. Опять линии погрешности являются прямыми.

Для $q > 1$ точный ответ неизвестен. В этом случае для получения значений погрешности можно воспользоваться следующими соображениями. При закономерности (4.9) разности значений U при возрастании

N на единицу также должны ложиться на прямую в полулогарифмическом масштабе (это напоминает апостериорную оценку погрешности по методу Ричардсона для квадратур со степенной сходимостью). На рис. 3 приведены графики таких разностей для $q = 2, r = 1$. Они также оказываются прямыми. Все это позволяет сделать эвристическое

Утверждение 2. При выполнении условий Утверждения 1 погрешность формулы трапеций экспоненциально зависит от числа узлов сетки N . \square

Попробуем выяснить, от чего зависит коэффициент β в (4.9). Он не должен зависеть от максимумов модулей каких-либо производных $u(x)$, поскольку они входят в суммы формул Эйлера-Маклорена (4.6) и приводят к степенной сходимости. Поэтому рассмотрим гипотезу о связи β с особыми точками подынтегрального выражения.

Если скобка в знаменателе подынтегрального выражения в (номер формулы теста) обращается в нуль, то подынтегральное выражение имеет полюс порядка q . Это происходит при

$$\frac{c^2 + 1}{2c} = \cos(x) = \frac{e^{ix} + e^{-ix}}{2}. \quad (4.10)$$

Это уравнение имеет два решения: $e^{ix} = c$ или $e^{ix} = 1/c$. Следовательно, имеется две цепочки полюсов кратности q в точках

$$x^* = 2\pi m \pm i \ln(c), -\infty \leq m \leq +\infty. \quad (4.11)$$

Из рис № видно, что наименьшее расстояние между каким-либо из полюсов и ближайшей к нему точкой отрезка интегрирования есть $\ln(c)$.

Рис.1. Погрешность квадратуры трапеций для (18) при $p = 0$ и $q = 1$. Цифры около линий – величины s . Рис.2. Погрешность квадратуры трапеций для (18) при $p = 2$ и $q = 1$. Цифры около линий – величины s . Рис.3. Погрешность квадратуры трапеций для (16) при $p = 1$ и $q = 2$. Цифры около линий – величины s .

Предварительный просмотр графиков показал, что наклон $\beta \ln(c)$. Для тщательного анализа на рис.4 показано отношение $\beta/\ln(c)$ в зависимости от для нескольких значений $q = 1, 2$ и $r = 0, 1, 2$. Видно, что для полюса первого порядка ($q = 1$) при $r = 0$ это отношение с высокой точностью не зависит от s и равно 2. При $r = 1$ это отношение равно 2 при $\ln(c) \approx 0$; при увеличении s это отношение несколько уменьшается, причем линия в пределах графика близка к прямой. При $r = 2$ линия также начинается в точке 2, но её наклон ещё немного увеличивается.

Для полюсов второго порядка ($q = 2$) линии для разных r начинаются не со значения 2, а с несколько меньшего значения 1.93. Наклоны линий с разными r также несколько больше, причем даже линия с $r = 0$ уже имеет наклон. Наименьшее отношение в пределах указанного графика 1.70.

Рис.4. Зависимость $\beta/\ln(c)$ от величины c для различных p и q . Это позволяет сделать

Утверждение 3. Наклон β в (4.9) с хорошей точностью пропорционален расстоянию от отрезка интегрирования до ближайшего полюса интегрируемой функции в комплексной плоскости. Коэффициент пропорциональности достигает 2 в случае полюса первого порядка, и несколько уменьшается с увеличением кратности полюса и расстояния до него. \square

Практические рекомендации. 1. При использовании квадратурных формул со степенной сходимостью удобно сгущать сетки по N последовательно вдвое. Это позволяет использовать обычную процедуру Ричардсона для получения априорной асимптотически точной оценки погрешности. Такое сгущение экономично, поскольку суммарный объем всех расчетов лишь вдвое превышает объем расчетов на последней сетке [4].

Для квадратур с экспоненциальной сходимостью (4.9) также можно пользоваться процедурой Ричардсона, если сгущать сетки не вдвое, а каждый раз увеличивая N на 1. При этом будет получаться асимптотически точная апостериорная оценка погрешности. Однако такое сгущение сеток экономически невыгодно, поскольку суммарный объем вычислений будет в $\sim N/2$ раз больше, чем расчет на последней сетке.

Поэтому в практических расчетах удобнее увеличивать N в 2 раза. Из (4.9) нетрудно получить, что при этом $\delta_{2N} \sim \delta_N^2$. Такой закон убывания напоминает сходимость ньютоновских итераций вблизи простого корня: число верных десятичных знаков приблизительно удваивается с увеличением N в 2 раза. Поэтому на практике останавливаются на такой сетке $2N$, когда отклонение от результата на предыдущей сетке становится меньше $\varepsilon^{2/3}$, где ε – ошибка единичного округления компьютера.

2. Для формулы трапеций на равномерной сетке полезен следующий прием, вдвое уменьшающий трудоемкость вычислений. На сетке с N узлами и шагом h формула трапеций имеет вид

$$U_N = h\left(\frac{u_0}{2} + u_1 + \dots + u_{(N-1)} + \frac{u_N}{2}\right). \quad (4.12)$$

При удвоении сетки все узлы предыдущей сетки становятся четными узлами новой сетки, и заново вычислять значения функций в них не надо. Достаточно найти значения функции в новых (нечетных) узлах и вычислить

$$U_{2N} = \frac{1}{2}U_N + \frac{h}{2}(u_1 + u_3 + u_5 \dots + u_{2N-1}). \quad (4.13)$$

где нечетные индексы относятся к узлам новой сетки.

§4.4. Сравнение теории с эвристическими оценками

Теорема 1 с оценкой погрешности (4.1) и эвристические Утверждения 1-3, полученные из численных экспериментов, качественно близки. Однако между ними имеется ряд различий. Обсудим их.

Теорема 1 строго доказана. В ней рассмотрен интеграл от периодической аналитической функции, взятой по полному периоду. Получена мажорантная оценка погрешности, справедливая для любого типа особых точек, включая существенно особые.

При получении эвристических Утверждений 1-3 также неявно предполагалась аналитичность функции. Однако периодичность функции не предполагалась. Правда, в рассмотренном тесте функция была периодической, но интеграл брался только по половине периода. Поэтому на самом деле требовалось лишь равенство нечетных производных на концах отрезка интегрирования (такое обобщение существенно для практических применений). Однако в численных экспериментах рассматривался только простейший случай особых точек - полюса первого и второго порядков. Это более благоприятная ситуация, позволившая получить более сильные оценки, причем асимптотически точные, а не мажорантные.

И Теорема 1, и Утверждения 1-3 дают экспоненциальную сходимость. Однако в Теореме 1 коэффициент перед числом интервалов N равен расстоянию до ближайшей особой точки. В эвристических оценках этот коэффициент почти в 2 раза больше, то есть реальная сходимость гораздо быстрее теоретической оценки (4.1). Однако, при увеличении порядка полюса, этот коэффициент уменьшается. Поэтому возможно, что для существенно особых точек он уменьшится до теоретического значения.

В знаменателе теоретической оценки (4.1) из экспоненты вычитается 1. В эвристических оценках этого не наблюдалось. Скорее всего, вычитание 1 связано с особенностями теоретического вывода. Заметим, что это вычитание существенно лишь, когда особая точка стремится к отрезку интегрирования. Если расстояние до особой точки не очень мало, то при разумном числе интервалов N экспонента будет существенно превышать 1 и указанный эффект станет незаметным.

§4.5. Квадратуры для функций ФД полуцелых индексов

Сравнительно экономичным способом прямого вычисления функций ФД могут служить квадратуры с экспоненциальной сходимостью. При этом возникает два различных случая, которые опишем ниже.

Функции индекса $k \geq -1/2$. Они определяются через сходящийся интеграл (1.5). Сделаем в интеграле (1) замену переменных $t = \tau^2$. Тогда интеграл (1) приведет к следующему виду

$$I_k(x) = 2 \int_0^\infty \frac{\tau^{2k+1} d\tau}{1 + \exp(\tau^2 - x)}, k \geq -\frac{1}{2}. \quad (4.14)$$

При полуцелых $k \geq -1/2$ показатель степени в подынтегральном выражении будет целым четным неотрицательным числом. Поэтому подынтегральное выражение будет четной функцией τ , и все его нечетные производные на нижнем пределе интегрирования $\tau = 0$ обращаются в нуль. Тем самым, на нижнем пределе интегрирования удовлетворяется условие Частного случая Утверждения 1.

На верхнем пределе интегрирования подынтегральное выражение убывает как $\exp(-\tau^2)$. При этом все производные быстро стремятся к нулю. Но применять формулу Эйлера-Маклорена на равномерной сетке (4.6) к бесконечному интервалу невозможно. Поэтому для численного интегрирования надо обрезать интеграл (4.14) и положить

$$I_k(x) \approx 2 \int_0^T \frac{\tau^{2k+1} d\tau}{1 + \exp(\tau^2 - x)}, k \geq -\frac{1}{2}. \quad (4.15)$$

Верхний предел T нужно выбирать так, чтобы во-первых, отброшенной частью интеграла можно было бы пренебречь, во-вторых, чтобы производные при T были бы настолько малы, чтобы их вклад в формулы Эйлера-Маклорена был пренебрежимо мал. Тогда на отрезке $0 \leq \tau \leq T$ формула Эйлера-Маклорена на равномерной сетке (15) обеспечит экспоненциальную сходимость.

Вообще говоря, оценка минимального T зависит от x и k . Очевидно, T возрастает при увеличении x или увеличении k . Нетрудно оценить отброшенную часть интеграла (4.15). Для не малых $x > 0$ она составляет примерно $T^{2k} \exp(x - T^2)$. Для получения относительной ошибки эту величину нужно сравнить с асимптотикой интеграла (1.5) $x^{k+1}/(k+1)$. Для получения относительной погрешности ε должно выполняться условие

$$TODO \quad (4.16)$$

целесообразно брать T "с запасом" для обеспечения надежности алгоритма. Но левая часть (32) быстро возрастает при увеличении T , так что на практике достаточно умеренного увеличения T .

Если определять T по значению x_{min} , то такая величина будет пригодна для всех $x = x_{min}$. Рассмотрим два крайних случая при 10^{-16} (double precision). Первый соответствует $k = -1/2, x_{min} = 39$ (см. Таблицу 4); он дает $T = 8.4$. Второй соответствует $k = 7/2, x_{min} = 29$; он дает $T > 8.5$. Для остальных функций полуцелых индексов $k \geq -1/2$ получаются практически такие же результаты. Поэтому для всех индексов $k \geq -1/2$ и для любых аргументов $x \leq x_{min}$ далее будем единообразно брать $T = 12$. Это создает необходимый запас надежности, включая обеспечение малости высоких производных на правой границе.

Подробнее обсудим сходимость квадратур. На рис. 5 изображена комплексная плоскость переменной интегрирования τ . Отрезок интегрирования выделен жирной линией – это вещественная положительная по-

луось. Показатель экспоненциальной сходимости определяется расстоянием от ближайшего полюса до промежутка интегрирования. Найдем положение полюсов. Это полюсы первого порядка, возникающие при обращении в нуль знаменателя подынтегрального выражения. Это происходит при $\exp(\tau^2 - x) = -1$. Отсюда получается цепочка полюсов:

$$\tau_m^2 = x + i\pi(1 + 2m), -\infty \leq m \leq +\infty, -\infty \leq x \leq +\infty. \quad (4.17)$$

Извлечение квадратного корня из правой части дает

$$\operatorname{Re}\tau_m = \sqrt{\frac{x + \sqrt{x^2 + \pi^2(1 + 2m)^2}}{2}}, \operatorname{Im}\tau_m = \sqrt{\frac{-x + \sqrt{x^2 + \pi^2(1 + 2m)^2}}{2}} \quad (4.18)$$

Ближайший к промежутку интегрирования по τ полюс соответствует $m = 0$. Его расстояние от промежутка интегрирования равно

$$\operatorname{Im}\tau_0 = \sqrt{\frac{-x + \sqrt{x^2 + \pi^2}}{2}} = \pi\sqrt{2(x + \sqrt{x^2 + \pi^2})}, -\infty \leq x \leq +\infty. \quad (4.19)$$

При $x < 0$ удобнее пользоваться первым выражением, а при $x > 0$ - вторым. Видно, что

$$\operatorname{Im}\tau_0 \rightarrow \sqrt{|x|}, x \rightarrow -\infty \operatorname{Im}\tau_0 = \sqrt{\pi/2}, x = 0 \operatorname{Im}\tau_0 \rightarrow \pi/(2\sqrt{x}), x \rightarrow +\infty \quad (4.20)$$

Поведение зависимости $\operatorname{Im}\tau(x)$ изображено на рисунке. Таким образом, при $x \rightarrow -\infty$ расстояние до ближайшего полюса быстро возрастает, и формула трапеций обеспечивает высокую точность при небольшом числе узлов. Напротив, при $x \rightarrow +\infty$ расстояние до ближайшего полюса быстро уменьшается, и сходимость будет существенно медленнее. Однако даже в этом случае число узлов остается умеренным благодаря экспоненциальной скорости сходимости.

Рис.5. Полюса подынтегрального выражения для функций ФД.

Численные расчеты подтверждают эти соображения. Интеграл (4.15) вычисляется на отрезке $[0, 12]$ по формуле трапеций (4.6). Вычисление проводится с автоматическим сгущением сетки до выхода на ошибки округления. Эта процедура аналогична тому, что изложено в Практических рекомендациях в Разделе 4.4. При этом окончательные сетки содержат $N = 48$ или 96 узлов при $x = 0$ и доходят до $N = 192$ при $x = 30 - 40$. Действительно видно, что число узлов слабо зависит от x при широких пределах изменения x .

Индекс $k = -3/2$. Для этого индекса интеграл (1) оказывается расходящимся. Поэтому для вычисления функций целесообразно переопределить ее через производную функции старшего индекса и провести сле-

дующее преобразование [9]

$$I_{-\frac{3}{2}}(x) = -2I'_{-\frac{1}{2}}(x) = -2 \int_0^{\infty} \frac{d}{dx} \left(\frac{1}{1+e^{t-x}} \right) \frac{dt}{\sqrt{t}} = -2 \int_0^{\infty} \frac{e^{t+x}}{(e^t + e^x)^2} \frac{dt}{\sqrt{t}} \quad (4.21)$$

В последнем выражении сделаем замену переменных $t = \tau^2$ и одновременно обрежем интеграл по верхнему пределу, аналогично показанному выше. Получим выражение

$$I_{-\frac{3}{2}} \approx -4 \int_0^T \frac{e^{\tau^2+x} d\tau}{(e^{\tau^2} + e^x)^2}. \quad (4.22)$$

Подынтегральное выражение есть четная функция τ , так что при достаточно большом T квадратура трапеций на равномерной сетке для интеграла (34) будет сходиться экспоненциально. Нетрудно видеть, что полюсы подынтегрального выражения также находятся согласно формуле (номер формулы для полюсов подынтегрального выражения), но теперь это будут полюсы второго порядка. Выше мы видели (см. рис № зависимости от кратности полюса), что скорость сходимости для полюсов второго порядка лишь незначительно медленнее, чем для полюсов первого порядка. Также можно взять $T = 12$, что обеспечивает достаточный запас надежности. Тогда расчеты по формуле трапеций для получения точности $\varepsilon = 10^{-16}$ требуют $N = 96$ узлов сетки при $x = 0$ и $N = 384$ при $x \approx 44$.

Практические рекомендации. В принципе можно подбирать T своим для каждого x : это уменьшает необходимое число интервалов N для небольших x . Однако на практике можно получить в среднем гораздо меньшую трудоемкость следующим образом.

При вычислении выражения самими трудоемкими операциями, являются вычисления экспонент. Величина e^x не зависит от узла сетки, и её следует вычислять только 1 раз. Величина $\exp(\tau_n^2)$ зависит от узла сетки, однако выбирая $T = 12$ для всех функций и всех сеток, мы обеспечиваем двойную точность (double precision) практически для всех интересующих нас значений x . Поэтому можно заранее разделить на $N = 384$ интервала, чего также достаточно практически для любых значений x . На этой сетке вычислим все значения $\exp(\tau_n^2)$, $0 \leq n \leq 384$, и включим эту заранее вычисленную таблицу в программу вычисления функций ФД. Тогда на любой конкретной сетке достаточно делать только выборку необходимых значений экспонент. Этот несложный прием многократно уменьшает время расчета.

Следует также использовать вариант вычисления по формуле трапеций, описанный в разделе 4.4. Это ещё в 2 раза уменьшает трудоемкость вычислений.

Напомним, что суммирование слагаемых в формуле трапеций надо начинать не с левого, а с правого конца, т.е. с меньших значений подынтегрального выражения. Это уменьшает ошибки округления. При больших N выигрыш в точности может оказаться существенным.

§4.6. Квадратуры для интегральной функций ФД

Интегральная функция ФД определяется как интеграл (1.6) от квадрата $I_{-1/2}(x)$. Даже если мы умеем хорошо вычислять функцию $I_{-1/2}(x)$, находить интегральную функцию прямой квадратурой (1.*) невыгодно. Такая квадратура будет иметь только степенную сходимость, и для получения точности $\varepsilon = 10^{-16}$ потребуется неприемлемо большое число узлов сетки. Однако, оказалось возможным свести задачу к экспоненциально сходящимся квадратурам.

Для этого подставим в (1) выражение $I_{-1/2}(\xi)$ через одномерный интеграл, сразу делая замену переменной интегрирования $t = \tau^2$. Квадрат такого последнего одномерного интеграла будет двойным интегралом по $d\tau d\theta$. Поэтому окончательно, вместо одномерного интеграла (1.) получим тройной интеграл

$$J(x) = 4 \int_0^x d\xi \int_0^\infty \int_0^\infty \frac{d\tau d\theta}{[1 + \exp(\tau^2 - \xi)][1 + \exp(\theta^2 - \xi)]} \quad (4.23)$$

Переменим порядок интегрирования, и сначала произведем интегрирование по ξ . Умножая числитель и знаменатель подынтегрального выражения на $e^{2\xi}$, преобразуем его к следующему виду:

$$\begin{aligned} \frac{d\xi}{[1 + \exp(\tau^2 - \xi)][1 + \exp(\theta^2 - \xi)]} &= \frac{e^\xi de^\xi}{(e^\xi + e^{\tau^2})(e^\xi + e^{\theta^2})} = \\ &= \frac{1}{e^{\tau^2} - e^{\theta^2}} \left(\frac{e^{\tau^2}}{e^{\tau^2} + e^\xi} - \frac{e^{\theta^2}}{e^{\theta^2} + e^\xi} \right) d(e^\xi). \end{aligned} \quad (4.24)$$

Теперь интегрирование по ξ выполняется в элементарных функциях, и тройной интеграл (4.23) превращается в двойной:

$$J(x) = 4 \int_0^\infty \int_0^\infty \frac{e^{\tau^2} \ln(1 + e^{x-\tau^2}) - e^{\theta^2} \ln(1 + e^{x-\theta^2})}{e^{\tau^2} - e^{\theta^2}} d\tau d\theta \quad (4.25)$$

Подынтегральная функция симметрична и положительна. При $\tau = \theta$ в ней возникает неопределенность типа $0/0$, которая раскрывается по правилу Лопиталя:

$$\left. \frac{e^{\tau^2} \ln(1 + e^{x-\tau^2}) - e^{\theta^2} \ln(1 + e^{x-\theta^2})}{e^{\tau^2} - e^{\theta^2}} \right|_{\tau=\theta} = \ln(1 + e^{x-\tau^2}) - \frac{e^x}{e^{\tau^2} + e^x}. \quad (4.26)$$

Вблизи линий $\tau = \theta$ при непосредственном вычислении подынтегрального выражения (4.25) на разностной сетке может возникать потеря точности; в этом случае в окрестности линии $\tau = \theta$ надо использовать уточнение выражения (??) с использованием более высоких производных числителя.

Подынтегральное выражение (4.25) четно по τ и по θ . Тем самым, на полуосях $\tau = 0$ и $\theta = 0$ выполняются условия Утверждения 2 для экспоненциальной сходимости квадратур трапеций. При $\tau^2 \gg x$ или $\theta^2 \gg x$ подынтегральное выражение очень быстро убывает со всеми производными. Поэтому можно ограничить область интегрирования квадратом $0 \leq \tau, \theta \leq T$. Следовательно, двумерная квадратура трапеций на равномерной сетке в этом квадрате будет иметь экспоненциальную сходимость. Поскольку подынтегральное выражение симметрично, можно ограничиться интегрированием по треугольнику $0 \leq \tau \leq \theta \leq T$, что вдвое уменьшает объем вычислений.

Обозначим подынтегральное выражение через $f(\tau, \theta)$ и введем равномерные сетки $h = T/N$, $\tau_n = nh$, $\theta_m = mh$. Тогда квадратура трапеций для треугольной области запишется с весами w_{nm} :

$$J(x) = 8h^2 \sum_{n=N}^0 \sum_{m=N}^n w_{nm} f(\tau_n, \theta_m); \quad (4.27)$$

$$w_{00} = \frac{1}{8}, w_{n0} = w_{nn} = \frac{1}{2}, n > 0; w_{nm} = 1, n > 1, n > m.$$

Вес на верхней границе треугольника безразличен, так как там функция и ее производные пренебрежимо малы. Суммирование в формуле (4.27) поставлено в обратном порядке, так как суммирование от малых членов к большим уменьшает ошибки округления.

Замечание. В разделе 4.5 отмечено, что вычисление экспонент является довольно трудоемкой операцией. Поэтому следует заранее вычислить значения экспонент $\exp(\tau^2)$, $0 \leq n \leq N$. Далее эти значения надо подставлять при вычислении подынтегрального выражения при соответствующих значениях $\tau_n \theta_m$. Это примерно в N раз уменьшает объем вычислений. Видно, что заранее вычисленные экспоненты одинаково работают для обоих направлений интегрирования. Поскольку остальные действия являются арифметическими, то трудоемкость вычисления двумерного интеграла будет несильно отличаться от трудоемкости вычисления одномерного интеграла. При таком усовершенствовании вычисление двумерного интеграла оказывается пригодным, как способ прямого вычисления функции во встроенных стандартных подпрограммах.

Для обеспечения точности $\varepsilon = 10^{-16}$ при $x = 0$ достаточно сетки $N = 96$, а вблизи x_{min} – сетки $N = 384$.

§4.7. Некоторые результаты.

Таблица 4.1: Реперные значения функций ФД при $x = 0$

k	$I_k(0)/\Gamma(k+1)$	k	$I_k(0)/\Gamma(k+1)$
-3/2	0.380104812609684	2	0.901542677369696
-1/2	0.604898643421630	5/2	0.927553577773948
0	0.693147180559945	3	0.947032829497246
1/2	0.765147024625408	7/2	0.961483656632978
1	0.822467033424113	4	0.972119770446909
3/2	0.867199889012184		

Приведем некоторые результаты численных расчетов. Заметим, что при $x \rightarrow -\infty$ главный член простейшего ряда есть $I_k(x)_{app} = \Gamma(k+1)e^x$. Поэтому отношение $I_k(x)/\Gamma(k+1)_{app} = e^x$, т.е. соответственно нормированная функция ФД при любых k практически совпадают для $x \rightarrow -\infty$. Фактически, эти отношения остаются довольно близкими даже при $x = 0$. В Таблице 7 приведены значения этих нормированных функций с 16-ю десятичными знаками при $x = 0$. Видно, что эти значения слабо возрастают с увеличением k . Приведенные значения полезны как реперные точки. Таблица 7

На Рис.6 приведены графики нормированных функций ФД. Поскольку все нормированные функции положительны и меняются в очень больших пределах, то для рисунка выбран полулогарифмический масштаб. В левой части графика при $x \rightarrow -\infty$ все нормированные функции быстро стремятся к прямой, проходящей через начало координат. В правой части они расходятся, не пересекаясь; чем больше k , тем выше лежит кривая. Для $k > -1$ каждая кривая является монотонно возрастающей; но для $k = -3/2$ кривая имеет максимум.

Глава 5

ФУНКЦИИ ПОЛУЦЕЛОГО ИНДЕКСА

Глава 6

АППРОКСИМАЦИИ ФУНКЦИЙ ФД

§6.1. Глобальные аппроксимации

Обобщенные формулы. Для несложных оценочных расчетов физикам полезно иметь простые формулы, непрерывно и гладко описывающие функции ФД во всем диапазоне значений аргумента $-\infty \leq x \leq +\infty$. Для успешного применения такие формулы должны описывать главные члены левых и правых асимптотик функций. Физически это обеспечивает переход фермионов в идеальный газ при высоких температурах, и в полностью вырожденный газ при низких температурах.

Для построения таких формул учтем следующие соображения. Левые асимптотики ($x \rightarrow -\infty$) функций ФД являются разложениями по e^x . Правые асимптотики ($x \rightarrow +\infty$) имеют главный член $\sim x^{k+1}$ и содержат разложения по степеням x^{-2} . Эти разложения сильно различаются по своей математической структуре. Левая асимптотика грубо искажает поведение функции ФД при $x > 0$, а правая при $x < 0$, поэтому построение оценочных формул в виде явной зависимости от x нетривиально.

Однако левые асимптотики всех функций ФД качественно сходны между собой, и правые асимптотики также качественно сходны между собой. Поэтому естественна идея выражать функции ФД индекса k через функции ФД индекса m . Такие формулы были предложены в [Ритус ЖВМ,препринт]. Рассмотрим их построение. При этом учтем, что когда аргумент x меняется от $-\infty$ до $+\infty$, значения функции ФД (при $k > -1$) монотонно возрастают от 0 до $+\infty$. При этом качественное поведение функций при $x \rightarrow +\infty$ и $x \rightarrow -\infty$ качественно различается.

При $x \rightarrow -\infty$ все функции ФД разлагаются в степенные ряды по e^x , причем главным членом является первая степень e^x (TODO ссылка на формулу). Это означает, что каждая функция разлагается в степенной

ряд по другой функции. Из () получаем следующее разложение:

$$I_k = \Gamma(k+1) \sum_{p=0}^{\infty} a_p \left[\frac{I_m}{\Gamma(m+1)} \right]^{1+p}. \quad (6.1)$$

Приведем главный и 2 следующих коэффициента этого разложения:

$$a_0 = 1, \text{TODO}. \quad (6.2)$$

Этот ряд сходится при достаточно малых I_m ; точное установление границ сходимости не производилось.

При $x \rightarrow +\infty$ главный член функции ФД имеет вид некоторой степени x . Он домножается на асимптотическое разложение по степеням x^{-2} . Тогда из (TODO ссылка нф) вытекает следующее разложение:

$$I_k \approx \frac{1}{k+1} \sum_{p=0} b_p [(m+1)I_m]^{\frac{k+1-2m}{m+1}}. \quad (6.3)$$

Его главный и 2 следующих коэффициента равны

$$b_0 = 1, \text{TODO}. \quad (6.4)$$

Это разложение является асимптотическим.

Построим следующую вполне гладкую аппроксимацию, выражающую одну функцию ФД через другую и правильно передающую главные члены и параметры малости обоих разложений (6.1) и (6.3):

$$\text{TODO} \quad (6.5)$$

Здесь первая сумма обеспечивает качественно правильное разложение при $I_m \rightarrow 0$, а вторая сумма делает то же при $I_m \rightarrow +\infty$. Заметим, что последний член первой суммы является её естественным окончанием. Но показатели степени во второй сумме подобраны так, что он может формально рассматриваться как член второй суммы. Это должно улучшить "сшивание" двух качественно отличающихся разложений.

Значения свободных параметров $c_p, 1 \leq p \leq N-1$, и числа членов в суммах можно варьировать для удовлетворения различных критериев. Например, соответствующим выбором младших коэффициентов $c_k, 1 \leq p \leq P_0 \leq P$, можно добиться точной передачи P_0 членов ряда (6.1), следующих за главным. Аналогично можно выбрать старшие коэффициенты $c_k, N-1 \geq p \geq P_1 \geq P$, так, чтобы правильно передать $N-P_1$ членов разложения (6.3). Остальные коэффициенты естественно определяются из условия минимума погрешности в $\|\cdot\|_C$ или $\|\cdot\|_{L_2}$, при этом для положительных и сильно меняющихся по величине функций ФД целесообразно минимизировать не абсолютную, а относительную погрешность.

Очевидно, увеличение N позволяет уменьшить погрешность аппроксимации, причем трудоемкость вычислений возрастает не сильно, ибо основное время расходуется на вычисление не сумм, а двух нецелых степеней. Значение P разумно полагать близким к $n/2$.

Замечания. 1° В задачах атомной физики аргумент x играет роль потенциала поля, в котором движется электрон, а функция $I_{\frac{1}{2}}(x)$ пропорциональна плотности электронов в соответствующей точке поля. Поэтому наибольшее практическое значение имеют формулы, выражающие функции ФД либо через $I_0(x)$ (в этом случае устанавливается зависимость от потенциала поля), как сделано выше, либо через $I_{\frac{1}{2}}(x)$, т.е. через плотность электронов.

2°. Формулу (6.5) затруднительно применять для получения прецизионных аппроксимаций. Прецизионные аппроксимации требуют большого числа коэффициентов. При этом задача определения методом наименьших квадратов или другим методом становится плохо обусловленной. Кроме того, при большом числе членов будет очень мал показатель степени k , что также ухудшает обусловленность задачи. Однако для построения гладких оценочных аппроксимаций невысокой точности с небольшим числом членов формула (6.5) удобна.

3°. Вычисление $I_k(x)$ через $I_m(x)$, где в свою очередь $I_m(x)$ непосредственно вычисляется по x , означает введение промежуточного аргумента. Наиболее удобным оказался выбор $m = 0$, потому что тогда промежуточный аргумент явно выражающаяся через элементарные функции от :

$$y(x) \equiv I_0(x) = \ln(1 + e^x); 0 \leq y < +\infty, -\infty < x < +\infty. \quad (6.6)$$

Такой аргумент ведет себя как e^x при $x \rightarrow -\infty$, и как x при $x \rightarrow +\infty$. Таким образом, он гладко склеивает экспоненту левой асимптотики с полиномом правой асимптотики. Поэтому целесообразно искать аппроксимацию $I_k(x)$ в виде некоторых явных зависимостей от промежуточного аргумента y .

Ниже построен большой набор таких приближенных формул.

Двучленные формулы. Двучленная формула была предложена в [8]. Запишем ее в несколько более удобной форме:

$$I_k(x) \approx \frac{y}{k+1} (\Gamma(k+2)^{1/k} + y)^k, 0 \leq y < +\infty. \quad (6.7)$$

Эта формула не содержит подгоночных параметров, её можно применять для любых индексов $k > -2$.

На рис. №. показаны погрешности двучленной формулы (6.7) для целых индексов k . Видно, что для $k > 0$ эта формула всюду завышает значение функции. Погрешность её значительна. Максимальное завышение составляет 30% для $k = 1$ и доходит до 190% для $k = 4$. Поэтому

двучленная формула пригодна только для очень грубых оценок.

Трехчленная формула была предложена в [9-10]. Ее также запишем в несколько удобной форме:

$$I_k(x) \approx \frac{y}{k+1} \left([\Gamma(k+2)]^{\frac{3}{k}} (1+cy) + y^3 \right)^{\frac{k}{3}}; \quad (6.8)$$

коэффициент c подбирается так, чтобы минимизировать относительную погрешность аппроксимации. Формула построена так, что при $y \rightarrow 0$ она точно передает главный член левой асимптотики и порядок малости следующего члена разложения (но не точный коэффициент в этом члене). При $y \rightarrow +\infty$ она точно передает главный член асимптотики и порядок малости следующего члена. Это улучшает качественное поведение аппроксимации.

Формула (6.8) оказалась несравненно лучше двучленной формулы. Наличие подгоночного параметра c позволяет сделать профиль погрешности знакопеременным. Подбирая из условия чебышевского альтернанса, мы минимизируем относительную погрешность. Соответствующие оптимальные значения приведены в табл.2; в ней приведены значения для целых индексов k , взятых из [TODO ссылка на уточнение прецизионных аппроксимаций], и полуцелых индексов k , взятых из [Ритус ЖВМб препринт] (для $k=-1/2$, $k=1/2$, $k=3/2$) и [master Диплом] ($k=3/2$, $5/2$, $7/2$). Видно, что коэффициенты монотонно убывают с ростом k . В табл.2 приведены также погрешности полученных формул в процентах; они тем больше, чем сильнее отличается k от нуля. Эти погрешности много меньше, чем для формулы (6.7), и даже для $k=4$ не превышают 6.4%. Видно, что трехчленную формулу (6.8) уже можно рекомендовать для удовлетворительных оценочных расчетов.

Типичный график относительной погрешности этой формулы для $k=2$ приведен на рис.№. Видно, что погрешность знакопеременна, имеет 2 равных по модулю экстремума, т.е. удовлетворяет условию чебышевского альтернанса, и стремится к нулю на бесконечности.

Обобщенные трехчленные формулы. Запишем частный случай обобщенных формул (ссылка на номер) с тремя членами, правильно пе-

Таблица 6.1: Коэффициенты и погрешности формулы (6.8).

k	c	$\Delta_{max}(\%)$
-3/2		
-1/2	1.62	0.7
1/2	1.18	0.8
1	1.01	1.6
3/2	0.87	2.6
2	0.77	3.2
5/2	0.67	4.0
3	0.60	4.8
7/2	0.54	5.6
4	0.48	6.4

ределяющие главные члены левой и правой асимптотик:

$$\begin{aligned}
I_k &= \alpha_0 I_m \cdot (1 + \alpha_1 I_m + \alpha_2 I_m^\eta)^\kappa, \\
\alpha_0 &= \frac{\Gamma(k+1)}{\Gamma(m+1)}, \\
\alpha_2 &= \left[\frac{\Gamma(m+1)}{\Gamma(k+1)} \frac{(m+1)^{\frac{k+1}{m+1}}}{k+1} \right]^{1/\kappa}, \\
\kappa &= \frac{k-m}{m+3}, \\
\eta &= \frac{(m+3)}{m+1}.
\end{aligned} \tag{6.9}$$

Коэффициент α_1 подбирается для минимизации относительной погрешности. Погрешность знакопеременна и имеет два экстремума, удовлетворяющих условию чебышевского альтернанса.

Возьмем $m = 1/2$, что означает выбор электронной плотности в качестве аргумента. В таблице № приведены значения всех коэффициентов и погрешностей формулы (6.9) для индексов $k = -1/2, 0, 3/2$. Этот набор индексов обеспечивает расчеты статистической модели атома Томаса-Ферми; в частности, значение $k = 3/2$ дает выражение давления или энергии однородного электронного газа через его плотность. Видно, что достигается превосходная точность аппроксимации 0.6%. Этого достаточно для многих приложений, ибо неточность физических моделей зачастую превосходит эту величину.

Заметим, что такая точность при единственном свободном параметре свидетельствует об удачном выборе вида аппроксимации.

Пятичленные формулы дают еще более хорошие результаты. Запишем их так, чтобы члены с первого по третий выглядели как разложение по степеням $y \equiv I_0$ при $y \rightarrow 0$, а члены с третьего по пятый – как разложение по степеням y^{-2} при $y \rightarrow +\infty$:

$$I_k(x) \approx \frac{y}{k+1} (\Gamma(k+2)^{\frac{6}{k}} (1 + c_1 y + c_2 y^2) + c_3 y^4 + y^6)^{\frac{k}{6}}. \quad (6.10)$$

Выбирая коэффициенты из различных соображений, рассмотрим три варианта этой формулы.

1°. Улучшенные асимптотики. Выберем коэффициенты c_1 и c_3 так, чтобы правильно передать вторые члены левой и правой асимптотик:

$$c_1 = 3(1 - 2^{-k})/k, c_3 = \pi^2(k+1); \quad (6.11)$$

коэффициент c_2 оставим в качестве подгоночного. Такая формула обеспечивает описание не только пределов идеального и вырожденного ферми-газа, но и ближайшей поправки при уменьшении идеальности либо снятии вырождения. Особенно важен коэффициент c_3 , поскольку он позволяет описать тепловые свойства почти вырожденного газа (например, теплоемкость или проводимость при малых температурах).

Один подгоночный коэффициент c_2 может обеспечить лишь один нуль погрешности. Подбираем c_2 из условия чебышевского альтернанса. Это позволяет минимизировать погрешность. В табл.3 приведены оптимальные значения коэффициентов c_2 для целых и полуцелых значений k , а также погрешности полученных формул в процентах. Полученные погрешности составляют не более 2% даже для больших k , что втрое лучше, чем для формулы (7). Тем самым эта формула предпочтительнее для оценочных расчетов.

2°. Низкотемпературная асимптотика. Передача второго члена асимптотики при высоких температурах обычно не столь важна. Поэтому можно коэффициенты c_1 и c_2 сделать свободными параметрами, а коэффициент c_3 сохранить согласно (6.11). Это позволяет ввести второй нуль в график погрешности, а коэффициенты c_1 и c_2 подобрать из условия чебышевского альтернанса. Значения этих коэффициентов и соответствующие им значения максимальных погрешностей приведены в табл.4. Видно, что при этом относительная точность не хуже 1%. Это вдвое лучше, чем в табл.3. Таблица 3. Коэффициенты и погрешности формулы (8), вариант а). Таблица 4. Коэффициенты и погрешности формулы (8), вариант б) .

3°. Наилучшая точность. Пусть важна минимальная погрешность, а вторым членом правой асимптотики можно пожертвовать. Тогда все три коэффициента c_1, c_2, c_3 можно использовать как подгоночные и выбирать из условия чебышевского альтернанса. График погрешности при этом будет иметь три нуля. Соответствующие значения коэффициентов и погрешностей приведены в табл.№. Видно, что погрешности еще

Таблица 6.2: Коэффициенты и погрешности формулы (6.11), вариант в).

k	c_1	c_2	c_3	$\delta_{max}(\%)$
-3/2				
-1/2	1.846	5.430	7.166	0.28
1/2	1.44	2.47	16.58	0.14
1	1.28	1.78	21.50	0.20
3/2	1.14	1.32	26.60	0.28
2	0.99	1.02	31.42	0.30
5/2	0.87	0.801	36.513	0.41
3	0.78	0.65	41.45	0.45
7/2	0.70	0.53	46.430	0.47
4	0.63	0.44	51.59	0.50

вдвое уменьшаются по сравнению с табл. 4 и не превышают 0.5%, а для функций с небольшими индексами составляют $\sim 0.2\%$. Это превосходная точность для столь простых формул.

Отметим, что подобранные значения z оказались близкими к теоретическим значениям (6.11). Во-первых, это свидетельствует об удачном выборе аппроксимации. Во-вторых, это означает, что полученными формулами можно удовлетворительно пользоваться для описания теплоемкости и других аналогичных свойств почти вырожденного электронного газа.

§6.2. Прецизионные аппроксимации для целых k при $x \leq 0$

Вид аппроксимации. Ранее говорилось, что значение функций целого индекса при $x > 0$ выражается через значение этой функции при $x < 0$ при помощи несложной формулы (6.11). Поэтому нужно иметь способы прецизионного вычисления $I_k(x)$ при целых k для $x \geq 0$.

Таковыми прецизионными формулами могут служить всюду сходящиеся ряды (6.12). Однако для получения относительной погрешности $\varepsilon = 10^{-16}$ при $x \sim 0$ требуется суммировать 40 членов ряда. Это может оказаться недостаточно экономичным для стандартных программ. Поэтому рассмотрим, как строить более экономичные прецизионные формулы.

Используем следующие соображения. При $x \rightarrow -\infty$ главный член асимптотики есть $I_k(x) \approx \Gamma(k+1)e^x$. При x асимптотика имеет не экспоненциальное, а степенное поведение: $I_k(x) \approx x^{k+1}/(k+1)$. Поэтому попробуем взять в качестве аргумента $y = I_0(x)$.

Аппроксимация многочленом редко бывает удачной, хотя она широко используется в литературе. Обычно лучшие результаты дает рациональная

нальная аппроксимация, то есть приближение отношением многочленов. Такая аппроксимация может передать разные асимптотики функций. В данном случае мы выбрали следующую аппроксимацию:

$$I_k(x) \approx \Gamma(k+1)y \left(\frac{\sum_{n=0}^{N+1} a_n y^n}{\sum_{n=0}^N b_n y^n} \right)^k, a_0 = 1, b_0 = 1, x \leq 0. \quad (6.12)$$

При $x \rightarrow -\infty$ аппроксимация (6.12) точно передает первый член ряда (4). Если провести разложение (6.12) по степеням e^x при $x \rightarrow -\infty$, то оно качественно будет подобно ряду (4). При $x \rightarrow +\infty$ главный член разложения (6.12) будет $I_k(x) \sim x^{k+1}$, хотя коэффициент при нем не будет совпадать с точным. Однако даже такая передача правой асимптоты улучшает точность аппроксимации.

Нахождение коэффициентов. Алгоритмы вычисления коэффициентов рациональных аппроксимаций обычно достаточно сложны, если добиваться минимизации некоторой нормы погрешности. Наилучшим был бы алгоритм, обеспечивающий чебышевский альтернанс для относительной погрешности. Однако неясно, как строить такой алгоритм. В данном случае мы ограничились несложным эвристическим алгоритмом, дающим хорошие результаты. Преобразуем (6.12) к следующей форме:

$$z = \left[\frac{I_k(x)}{\Gamma(k+1)y} \right]^{\frac{1}{k}}, z \approx \left(\frac{\sum_{n=0}^{N+1} a_n y^n}{\sum_{n=0}^N b_n y^n} \right). \quad (6.13)$$

Аппроксимация (6.13) содержит $2N+1$ свободный коэффициент: a_n , с $1 \leq n \leq N+1$ и b_m , с $1 \leq m \leq N$. Мы ищем аппроксимацию на интервале $-\infty < x \leq 0$, то есть $0 < y \leq y_{\max} = \ln 2$. Выберем некоторым образом $2N+1$ узлов $y_j : 0 < y_1 < y_2 < \dots < y_{2N+1} = y_{\max}$. По величинам y_j вычислим соответствующие значения $x_j, I_k(x_j)$ и w_j . Потребуем, чтобы в j -х точках приближенное равенство (6.13) становилось точным, то есть поставим задачу интерполяции по выбранным узлам.

Заметим, что можно формально взять $j=0$ и положить $y_0=0$. Но фактически делать этого не нужно. В этой точке приближенное равенство (6.13) автоматически становится точным, поскольку аппроксимация (6.12) точно передает главный член левой асимптотики.

Система линейных уравнений для нахождения свободных коэф-

фициентов имеет следующий вид:

$$\sum_{j=1}^{2N+1} AY_j = B, \quad (6.14)$$

где

$$TODO \quad (6.15)$$

Обусловленность системы быстро ухудшается с ростом N . Тем не менее, существующая программа дает разумные результаты при её решении. Это может повлиять на величины вычисляемых коэффициентов, однако слабо влияет на погрешность полученных формул (это известный парадокс метода наименьших квадратов).

Узлы интерполяции. Положение выбранных узлов сильно влияет на качество полученной интерполяции. При неудачном положении этих точек могут возникать отрицательные коэффициенты a_n, b_m ; это опасно, особенно если знаменатель или числитель обращаются в нуль внутри требуемого диапазона значений y . При наличии таких полученная аппроксимация совершенно неприемлема. Мы опробовали на практике некоторые способы выбора точек интерполяции. Проиллюстрируем их на примере $k = 2$ и $N = 3$ (это означает 7 свободных коэффициентов); результаты для других k и N были аналогичными.

1°. **Линейное распределение.** Простейшим способом было линейное расположение узлов:

$$y_j = \frac{j \cdot y_{max}}{2N + 1}, 0 \leq j \leq 2N + 1; \quad (6.16)$$

напомним, что формально мы можем вводить y_0 , хотя в расчетах оно не используется. Профиль относительной погрешности δ для этого случая приведен на рис. №. Видно, что погрешность обращается в нуль во всех узлах интерполяции, а между ними имеет вид полуволн. Амплитуды этих полуволн невелики в средних интервалах, а в крайних они в несколько раз больше. Расчеты показывают, что при увеличении N отношение погрешности в центре и на периферии быстро возрастает. Это говорит о том, что в середине отрезка следовало бы увеличить расстояние между узлами интерполяции, а на краях отрезка уменьшить.

2°. **Чебышевское распределение.** Хорошо разработана теория аппроксимаций многочленами, наилучшая в норме C . В ней положение узлов интерполяции точно не вычисляется. Однако оно близко к распределению, описываемому тригонометрической функцией:

$$y_j = \frac{1}{2} y_{max} \left(1 - \cos \left(\frac{\pi j}{2N + 1} \right) \right), 0 \leq j \leq 2N + 1. \quad (6.17)$$

Профиль погрешности δ для выбора узлов (6.17) так же показан на рис.2. Между узлами интерполяции он имеет вид полуволн, амплитуды которых велики в середине отрезка и малы по краям. Поэтому для распределения (6.17) надо сблизить узлы интерполяции в середине отрезка и раздвинуть вблизи границ отрезка.

Заметим, что это не противоречит теоретическим результатам для чебышевских приближений: они относятся к аппроксимации многочленами, а мы используем аппроксимацию рациональными функциями.

3°. Смешанное распределение. Представляется естественным построить распределение узлов интерполяции, промежуточное между линейным и тригонометрическим. Такая задача уже возникала в сверхбыстром итерационном методе решения систем эллиптических уравнений на прямоугольной сетке. Воспользуемся предложенным в [ссылка] линейно-тригонометрическим распределением:

$$y_j = \frac{y_{max}}{2} \left[\frac{2\gamma j}{2N+1} + (1-\gamma) \left(1 - \cos \left(\frac{\pi j}{2N+1} \right) \right) \right], 0 \leq j \leq 2N+1. \quad (6.18)$$

Чему равно γ ? Пусть есть производящая линейная функция $F = \gamma F_1 + (1-\gamma)F_2$, где $F_1 = 2\zeta$, $F_2 = 1 - \cos(\pi\zeta)$, $\zeta \in [0; 1] \rightarrow F_2 \in [0; 2]$. Определим производную функции F :

$$\frac{dF}{d\zeta} = 2\gamma + \pi(1-\gamma) \sin(\pi\zeta). \quad (6.19)$$

При $\zeta = 0$, $\frac{dF}{d\zeta} = 2\gamma$, при $\zeta = 0.5$, $\frac{dF}{d\zeta} = 2\gamma + \pi(1-\gamma)$. Потребуем, чтобы производная в центре была бы вдвое больше, чем производная на границах, то есть $2 \frac{dF}{d\zeta} \Big|_{\zeta=0} = \frac{dF}{d\zeta} \Big|_{\zeta=0.5}$. Это дает $\gamma = \frac{\pi}{2+\pi}$.

Узлы (6.18) были построены для функции, которая фактически является отношением многочленов одинаковой степени. Поскольку у нас степени многочленов в числителе и знаменателе отличаются всего на 1, можно ожидать их хорошей применимости в нашем случае. Результаты расчета также показаны на рис.№. Видно, что теперь экстремумы в центральных и краевых интервалах почти одинаковы: отличие составляет $\sim 15\%$. Таким образом получаемое решение близко к чебышевскому альтернансу. Поэтому эвристическое распределение (6.18) можно считать почти неулучшаемым, и пользоваться им для любых N и k .

Видно также, что погрешность для линейно-тригонометрического распределения меньше в 7.5 раз, чем для линейного, и в 4 раза по сравнению с тригонометрическим. Это достаточно существенный выигрыш в точности.

Влияние числа параметров. Подробней исследуем погрешность аппроксимации для линейно- тригонометрических узлов. При $N \leq 3$

профили погрешности для всех k и N имеют тот же качественный вид, как и жирная линия на рис.2. Погрешность обращается в нуль в узлах интерполяции, между ними содержит $2N + 1$ полуволн, а экстремумы всех полуволн примерно одинаковы по модулю. Это показывает, что ошибки округления практически не влияют на погрешность аппроксимации.

Картина меняется при $N = 4$. Профиль относительной погрешности становится хаотическим с амплитудой $\sim 2 \cdot 10^{-16}$ (см. рис.3). Это показывает, что расчет вышел на ошибки округления, и дальнейшее увеличение числа параметров бессмысленно.

В Приложении приведены графики погрешностей для всех значений $k = 1 - 3$ и $N = 1 - 4$, используемых в данной работе. Они подтверждают сделанные выше выводы.

Погрешность аппроксимации можно характеризовать нормой $C : \delta_C = \max \|\delta(y)\|$. В табл.2 приведены значения этой величины (точнее, $lg(\delta_C)$) для различных k и N). Видно, что максимальная погрешность слабо зависит от k , но быстро убывает с увеличением N . Для $N = 1$ аппроксимация дает ~ 6 верных десятичных знаков, для $N = 2$ это ~ 10 знаков, для $N = 3$ это ~ 14 знаков; для $N = 4$ следовало бы ожидать ~ 18 знаков, но ошибки округления позволяют выйти всего лишь на ~ 16 знаков.