

4.2 Численные методы безусловной минимизации функции многих переменных

Ставится задача минимизации функции $f(\mathbf{x}) = f(x_1, x_2, \dots, x_n)$ в некоторой замкнутой области. Пусть $\mathbf{a} = (a_1, a_2, \dots, a_n)$ точка минимума $f(\mathbf{x})$. Будем говорить, что с точностью до ε точка \mathbf{x} может быть взята в качестве приближенного значения точки минимума, если

$$\|\mathbf{a} - \mathbf{x}\| < \varepsilon, \text{ где } \|\mathbf{a} - \mathbf{x}\| = \sqrt{\sum_{i=1}^n (a_i - x_i)^2} \text{ или} \\ \|\mathbf{a} - \mathbf{x}\| = \max |a_i - x_i| \quad (i = \overline{1, n}).$$

Для геометрической иллюстрации методов будем использовать функцию двух переменных. Напомним, что линиями уровня функции $Z = f(x, y)$ называют множество точек (x, y) , удовлетворяющих уравнению $f(x, y) = c$. Меняя c , мы получаем различные линии уровня функции $f(x, y)$. Геометрически линия уровня – это проекция на плоскость XOY линии пересечения $Z = f(x, y)$ и плоскости $Z = c$. Имея множество линий уровня, мы получаем представление о поведении $Z = f(x, y)$, говорят о рельефе функции $Z = f(x, y)$.

4.2.1 Методы многомерного прямого поиска

Суть методов многомерного прямого поиска, изложенных ниже, в том, что выбирают некоторую точку $\mathbf{x}^{(0)} = (x_1^{(0)}, \dots, x_n^{(0)})$ и допустимое направление поиска \mathbf{d} . Затем, отправляясь от точки $\mathbf{x}^{(0)}$ в направлении \mathbf{d} , минимизируют функцию одного переменного λ : $f(\mathbf{x} + \lambda \mathbf{d})$, изложенными выше методами.

Найдя $\lambda^{(0)}$, при котором $f(\mathbf{x}^{(0)} + \lambda \mathbf{d})$ получает минимальное значение, мы тем самым нашли точку $\mathbf{x}^{(1)} = \mathbf{x}^{(0)} + \lambda^{(0)} \mathbf{d}$, в которой значение f , вообще говоря, меньше чем в точке $\mathbf{x}^{(0)}$. Далее, отправляясь от $\mathbf{x}^{(1)}$ в некотором новом направлении \mathbf{d}_1 , получаем некоторую точку $\mathbf{x}^{(2)}$, в которой значение f вообще говоря, меньше чем

в $\mathbf{x}^{(1)}$ и т.д. Возникают вопросы: 1) Как целесообразнее выбирать \mathbf{d}_i ? 2) Сходится ли последовательность $\mathbf{x}^{(0)}, \mathbf{x}^{(1)}, \dots$? 3) Как оценить погрешность?

4.2.2 Метод циклического покоординатного спуска

Опишем алгоритм одного из таких методов, выбирая в качестве направлений поиска координатные векторы $\mathbf{d}_1 = \mathbf{l}_1, \mathbf{d}_2 = \mathbf{l}_2, \dots, \mathbf{d}_n = \mathbf{l}_n$, т.е. $\mathbf{d}_i = \mathbf{l}_i$ - вектор, все компоненты которого, за исключением i равны нулю.

Тогда

$$\mathbf{x}^{(0)} + \lambda \mathbf{l}_i = (x_1^{(0)}, x_2^{(0)}, \dots, x_{i-1}^{(0)}, x_i^{(0)} + \lambda, x_{i+1}^{(0)}, \dots, x_n^{(0)}),$$

$$f(\mathbf{x}^{(0)} + \lambda \mathbf{l}_i) = f(x_1^{(0)}, x_2^{(0)}, \dots, x_{i-1}^{(0)}, x_i^{(0)} + \lambda, x_{i+1}^{(0)}, \dots, x_n^{(0)}) \quad \text{и,}$$

следовательно, при минимизации функции $f(\mathbf{x}^{(0)} + \lambda \mathbf{l}_i)$, речь идет о минимизации функции одной переменной $x_i = x_i^{(0)} + \lambda_i$, при фиксированных остальных переменных.

Выберем начальную точку $\mathbf{y}_1 = \mathbf{x}^{(0)}$, направление $\mathbf{d}_1 = \mathbf{l}_1$ и, минимизируя $f(\mathbf{y}_1 + \lambda \mathbf{l}_1)$, найдем, что минимум этой функции достигается при λ_1 , в точке $\mathbf{y}_2 = \mathbf{y}_1 + \lambda_1 \mathbf{l}_1 = (x_1^{(0)} + \lambda_1, x_2^{(0)}, \dots, x_n^{(0)})$. Положим $\mathbf{y}_2 = \mathbf{y}_1 + \lambda_1 \mathbf{l}_1$, выберем направление \mathbf{l}_2 и, минимизируя $f(\mathbf{y}_2 + \lambda \mathbf{l}_2)$ найдем, что минимум этой функции достигается при λ_2 в точке $\mathbf{y}_3 = \mathbf{y}_2 + \lambda_2 \mathbf{l}_2 = (x_1^{(0)} + \lambda_1, x_2^{(0)} + \lambda_2, x_3^{(0)}, \dots, x_n^{(0)})$.

После n шагов найдем

$$\mathbf{y}_{n+1} = \mathbf{y}_n + \lambda_n \mathbf{l}_n = (x_1^{(0)} + \lambda_1, x_2^{(0)} + \lambda_2, \dots, x_n^{(0)} + \lambda_n).$$

Положим $\mathbf{x}^{(1)} = \mathbf{y}_{n+1}$ и тем самым завершим один цикл покоординатного спуска.

Отправляясь от $\mathbf{x}^{(1)}$, как от начальной точки, можем найти $\mathbf{x}^{(2)}$ и т.д.

Процесс завершить, если

$$\| \mathbf{x}^{(k+1)} - \mathbf{x}^{(k)} \| < \varepsilon,$$

где ε – требуемая точность,

$$\| \mathbf{x}^{(k+1)} - \mathbf{x}^{(k)} \| = \sqrt{\sum_{i=1}^n (x_i^{(k+1)} - x_i^{(k)})^2} \quad \text{или}$$

$$\| \mathbf{x}^{(k+1)} - \mathbf{x}^{(k)} \| = \max |x_i^{(k+1)} - x_i^{(k)}| \quad (i = \overline{1, n}).$$

Проиллюстрируем метод с помощью функции двух переменных (рисунок 4.4).

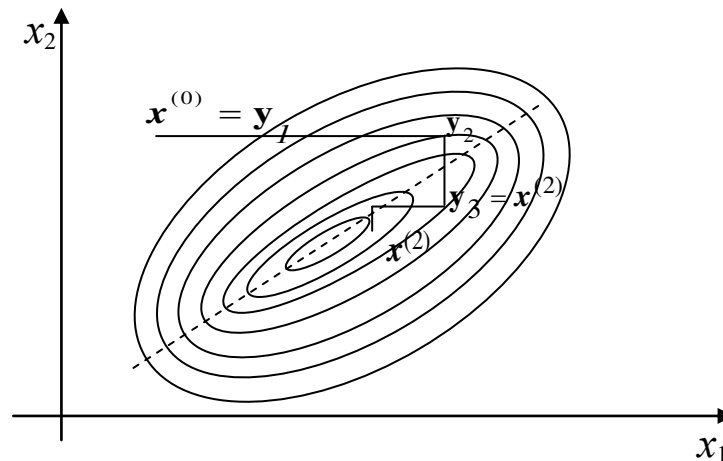


Рисунок 4.4 – Метод циклического покоординатного спуска

Двигаясь от точки $\mathbf{x}^{(0)} = \mathbf{y}_1$ параллельно оси Ox_1 , мы переходим от линий уровня с большим значением $f(x_1, x_2) = c$ к меньшим. Точка \mathbf{y}_2 , в которой мы коснемся некоторой линии уровня, является точкой минимума функции $f(x_1, x_2^{(0)})$ в направлении, параллельном оси Ox_1 . Двигаясь от точки \mathbf{y}_2 параллельно оси Ox_2 , мы, переходя от линий уровня с большим "с" к линиям уровня с меньшим "с", достигнем минимального значения в точке \mathbf{y}_3 - точке касания некоторой линии уровня. В точке \mathbf{y}_3 завершается цикл покоординатного спуска и получаем $\mathbf{x}^{(1)} = \mathbf{y}_3$.

Отправляясь от $\mathbf{x}^{(1)}$, как от начальной точки, можем найти $\mathbf{x}^{(2)}$ и т.д. Остается указать условия сходимости последовательности $\mathbf{x}^{(0)}, \mathbf{x}^{(1)}, \dots$ и указать оценку погрешности.

Для сходимости метода циклического покоординатного спуска достаточно следующих требований:

1) минимум $f(\mathbf{x})$ вдоль любого направления единственен;

2) последовательность точек $\mathbf{x}^{(0)}, \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(k)}$ принадлежит некоторому замкнутому ограниченному подмножеству области D .

Что касается погрешности, то ее можно определить по формуле: $\|\mathbf{a} - \mathbf{x}^{(k+1)}\| \leq \|\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}\|$.

ЗАМЕЧАНИЕ. Если функция f не является дифференцируемой в некоторых точках, то метод может остановиться в неоптимальной точке (рисунок 4.5).

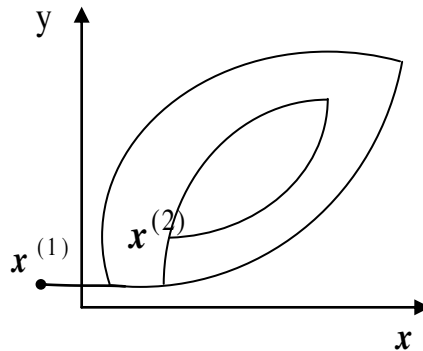


Рисунок 4.5 – Неоптимальная точка

Ясно, что поиск вдоль любой координатной оси в точке $\mathbf{x}^{(2)}$ не приводит к улучшению целевой функции. Это вызвано наличием так называемого оврага, то есть точки недифференцируемой функции $f(\mathbf{x})$. Разрешить эту ситуацию можно используя, например, метод Хука–Дживса.

4.2.3 Метод Хука-Дживса

Метод Хука-Дживса осуществляет два типа поиска: исследующий поиск и поиск по образцу. Выбрав начальную точку \mathbf{x}_1 методом циклического покоординатного спуска, находим точку \mathbf{x}_2 , т.е. осуществляем исследующий поиск. Если $\|\mathbf{x}_2 - \mathbf{x}_1\| < \varepsilon$, то точка минимума найдена, иначе осуществляем поиск по образцу в направлении $\mathbf{x}_2 - \mathbf{x}_1 = \mathbf{d}$,

что приводит нас в некоторую точку y . Приняв y за x_1 вновь проводим исследующий поиск и т.д. Сходимость метода Хука-Дживса обеспечена при тех же условиях, что и покоординатного спуска. Иллюстрация метода на рисунке 4.6 .

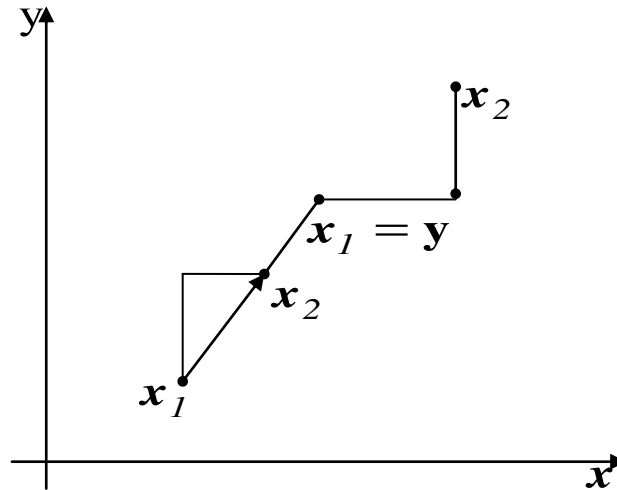


Рисунок 4.6 – Метод Хука-Дживса

4.2.4 Метод наискорейшего спуска

Известно, что направлением наибольшего возрастания функции f в точке $x^{(0)}$ является направление, задаваемое вектором $grad f(x^{(0)})$, а $(-grad f(x^{(0)}))$ задает направление наибольшего убывания функции f в точке $x^{(0)}$. Учитывая это, следует осуществлять линейный поиск не в направлении осей координат, а в направлении $-grad f$.

Рассмотрим алгоритм метода.

1. Выбираем x_i , $i=1$.
2. Если $\|grad f(x_i)\| < \varepsilon$, то x_i - искомая точка.
3. В противном случае положим $d_i = -grad f(x_i)$ и решив задачу о минимуме $f(x_i + \lambda d_i)$, найдем $\lambda_i > 0$.
4. Найдем $x_{i+1} = x_i + \lambda_i d_i$ и переходим к пункту 2, положив $i=i+1$.

Сходимость метода обеспечена, если f непрерывно дифференцируема, а генерируемая последовательность, принадлежит замкнутому ограниченному множеству.

Недостатком метода наискорейшего спуска является зачастую медленная сходимость в окрестности стационарной точки. Это очевидно, например, в тех случаях, когда линии уровня вытянуты в окрестности оптимальной точки. О функциях, поверхности уровня которых сильно вытянуты, говорят, что она имеет "овражный" характер. Геометрически медленная сходимость объясняется так: спустившись на "дно оврага" мы, двигаясь в направлении $(-\text{grad } f)$, будем переходить с одного склона оврага на другой, то есть зигзагообразно продвигаться к точке минимума (рисунок 4.7).

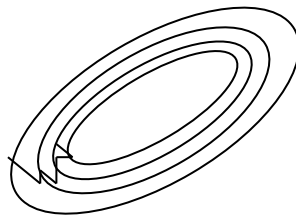


Рисунок 4.7 – Овражная функция

Что касается степени "овражности", то ее можно охарактеризовать с помощью минимального (λ_{min}) и максимального (λ_{max}) собственных чисел матрицы Гессе – $\left(\frac{\partial^2 f}{\partial x_i \partial x_j} \right) (i = \overline{1, n}; j = \overline{1, n})$. Чем меньше $\lambda_{min} / \lambda_{max}$, тем больше овражность. Более предпочтительными в таких случаях являются методы Ньютона и сопряженных направлений.

Коротко суть первого из них состоит в том, что функция $f(x)$ аппроксимируется (с помощью формулы Тейлора) многочленами второй степени, для которых находятся точки минимума. Последовательность таких точек приводит при определенных условиях к искомой точке минимума $f(x)$. Неудобством метода является необходимость многократного обращения матрицы Гессе.

Второй из методов для получения очередной точки требует проведения последовательной минимизации по каждому из n , специальным образом построенных направлений, которые называют сопряженными направлениями. Для квадратичной функции минимизация

вдоль " n " таких направлений позволяет "точно" достичь точки минимума, а следовательно можно ожидать хороших результатов и для достаточно гладких функций. Подробнее с этими методами можно познакомиться в [4].