

**Рейзлин В. И.,
доцент каф. ИПС, 414-КЦ**

**Методы оптимизации,
Экзамен.**

КОЛИЧЕСТВО КРЕДИТОВ

3 ECTS

Лекции

16 часов

Лабораторные занятия

32 часа

САМОСТОЯТЕЛЬНАЯ РАБОТА

60 часов

Итого

108 часов

В семестре за каждую посещенную лекцию начисляется 1 балл.

При выполнении **всех работ** и наборе 33 баллов студент допускается к экзамену.

Оптимизация в широком смысле слова находит применение в науке, технике и в любой другой области человеческой деятельности.

Оптимизация – целенаправленная деятельность, заключающаяся в получении наилучших результатов при соответствующих условиях.

Критерии оптимальности

Оптимизируемый вариант работы объекта должен оцениваться какой-то количественной мерой – *критерием оптимальности*.

Критерием оптимальности называется количественная оценка оптимизируемого качества объекта.

На основании выбранного критерия оптимальности составляется целевая функция, представляющая собой зависимость критерия оптимальности от параметров, влияющих на ее значение.

Вид критерия оптимальности или целевой функции определяется конкретной задачей оптимизации.

Таким образом, задача оптимизации сводится к нахождению экстремума целевой функции.

Любой оптимизируемый объект схематично можно представить в соответствии с рис. 1.

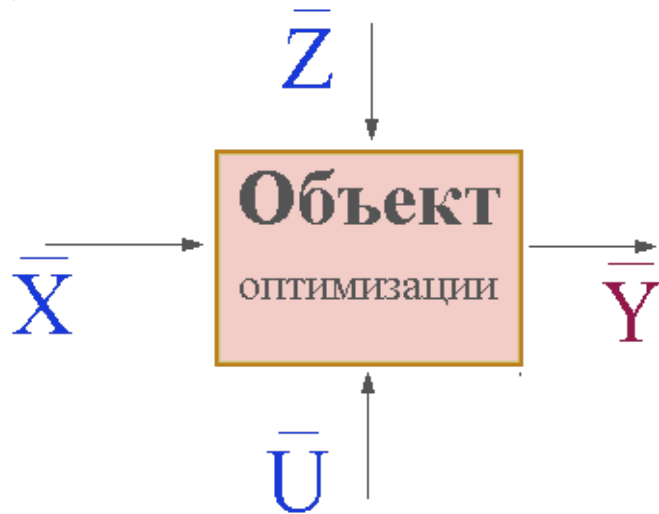


Рис.1. Оптимизируемый объект

\bar{Y} – выходы объекта;
 \bar{X} – контролируемые входные параметры;
 \bar{U} – регулируемые входные параметры
(управляемые параметры);
 \bar{Z} – неконтролируемые воздействия;

В том случае, когда случайные возмущения невелики и их воздействие на объект можно не учитывать, критерий оптимальности может быть представлен как функция *входных, выходных и управляемых* параметров:

$$R = R(X_1, X_2, \dots, Y_1, Y_2, \dots, U_1, U_2, \dots).$$

Так как $\bar{Y} = F(\bar{U})$, то при фиксированных \bar{X} можно записать:

$$R = R(\bar{U}).$$

При этом всякое изменение значений управляющих параметров двояко сказывается на величине R :

- прямо, так как управляющие параметры непосредственно входят в выражение критерия оптимизации;
- косвенно — через изменение выходных параметров процесса, которые зависят от управляющих.

Процедура решения задачи оптимизации обязательно включает, помимо выбора управляемых параметров, еще и установление ограничений на эти параметры (термостойкость, взрывобезопасность, мощность перекачивающих устройств и т.п.).

Итак, для решения задачи оптимизации необходимо:

а) выбрать критерий оптимальности для исследуемого объекта;

б) выбрать управляющие параметры;

в) установить ограничения на эти параметры;

г) построить модель объекта

$$Y = F(X, U),$$

которая устанавливает зависимости критерия оптимальности от всех аргументов и включает сопутствующие задаче ограничения;

д) на основе модели и критерия оптимальности определить целевую функцию

$$R = f(Y);$$

е) выбрать метод оптимизации, который позволит найти экстремальные значения ИСКОМЫХ ВЕЛИЧИН.

Оптимизация:

Дискретная:

Целочисленное программирование;
Стохастическое программирование.

Непрерывная:

Безусловная:

Глобальная
оптимизация;
Дифференцируемая
оптимизация;
Недифференцируемая
оптимизация.

Условная:

Линейное
программирование;
Нелинейные задачи;
Стохастическое
программирование.

Пример 1.

При протекании химических реакций для расчета изменения концентраций реагирующих веществ составляется кинетическая модель.

Например, для обратимой реакции

вещество A превращается в

вещество B



кинетическая модель представляется системой обыкновенных дифференциальных уравнений

$$\begin{cases} \frac{dC_A}{dt} = -k_1 C_A + k_2 C_B, \\ \frac{dC_B}{dt} = k_1 C_A - k_2 C_B, \end{cases} \quad (\text{B1})$$

с начальными условиями

$$t = 0, C_A = C_{A_0}, C_B = C_{B_0}.$$

Здесь C_A , C_B – концентрации веществ A и B , k_1 и k_2 – константы скоростей прямой и обратной реакций.

Различают *прямую* кинетическую задачу, когда константы k_1 и k_2 известны, и *обратную* кинетическую задачу, когда эти константы неизвестны и должны быть определены на основании экспериментальных данных.

Решение обратной кинетической задачи можно свести к оптимизационной процедуре следующим образом:

Пусть на отрезке времени $0 < t < T$ в точках $t_0 = 0, t_1, \dots, t_n = T$ известны экспериментальные данные для концентраций $C_{A_i}^{\text{Э}}, C_{B_i}^{\text{Э}}, i = 0, \dots, n$.

Решая численным методом прямую задачу для дифференциальной модели В1 при некоторых значениях констант k_1, k_2 , можно получить теоретические значения концентраций $C_{A_i}, C_{B_i}, i = 0, \dots, n$.

Можно найти значения констант k_1, k_2 , для которых теоретические значения для концентраций $C_{A_i}, C_{B_i}, i = 0, \dots, n$ будут отличаться от экспериментальных в некотором смысле меньше всего.

Например, следуя идеологии метода наименьших квадратов, можно построить целевую функцию

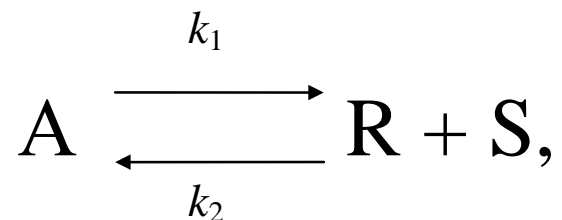
$$f(k_1, k_2) = \sum_{i=0}^n \left(C_{A_i} - C_{A_i}^{\text{Э}} \right)^2 + \sum_{i=0}^n \left(C_{B_i} - C_{B_i}^{\text{Э}} \right)^2$$

и найти ее минимум, т.е. искомые значения констант, каким либо численным методом оптимизации.

Пример 2.

Оптимизация работы каталитического химического реактора.

Если в реакторе протекает обратимая реакция



где R – целевой продукт, S – побочный продукт, то математическая модель реактора может быть описана системой уравнений материального и теплового балансов и имеет следующий вид:

$$\frac{dC_A}{dV} = \frac{(1-\varepsilon)}{\nu} \eta \left(-k_1 C_A + k_2 C_R C_S \right); \quad (\text{B.2})$$

при $V=0$; $C_A(0) = C_{A,0}$;

$$\begin{aligned}
\frac{dC_R}{dV} &= \frac{(1-\varepsilon)}{\nu} \eta \cdot (k_1 C_A - k_2 C_R \cdot C_S); \\
\frac{dC_S}{dV} &= \frac{(1-\varepsilon)}{\nu} \eta \cdot (k_1 C_A - k_2 C_R \cdot C_S); \\
\frac{dT}{dV} &= \frac{(-\Delta H) \cdot (1-\varepsilon) W \cdot \left(\frac{R_1 \cdot T}{P} \right)}{\nu \cdot C_P^{CM}}. \quad (\text{B.3})
\end{aligned}$$

Здесь V – объем слоя катализатора, м^3

ν – расход сырья, $\text{м}^3/\text{с}$;

W – наблюдаемая скорость химической реакции, $\text{кмоль}/\text{м}^3 \cdot \text{с}$;

$$W = \eta \cdot W_{\text{кин}};$$

$W_{\text{кин}}$ – скорость химической реакции в кинетической области, $\text{кмоль}/\text{м}^3 \cdot \text{с}$;

η – фактор эффективности;

T – температура в реакторе, К ;

C_p^{cm} – мольная теплоемкость потока, Дж/моль·К;

ε – порозность слоя катализатора;

$(-\Delta H)$ – тепловой эффект химической реакции,
Дж/моль;

P – давление в системе, МПа;

R_1 – газовая постоянная,
 $0,00845 \text{ м}^3 \cdot \text{МПа} / \text{кмоль} \cdot \text{К}$;

R – газовая постоянная, Дж/моль·К;

$$k_i = k_{i,0} \cdot e^{-E_i/RT}, \quad (\text{B.4})$$

где k_i – константа скорости i -й химической реакции, с^{-1} ;

E_i – энергия активации i -й химической реакции, Дж/моль.

$$C_P^{CM} = \sum_{i=1}^n C_{P_i} \cdot x_i, \quad (\text{B.5})$$

где x_i – мольная доля i -го компонента;

C_i - концентрация i -го вещества, кмоль/м³;

$$C_{P_i} = a_i + b_i T + C_i T^2, \quad (\text{B.6})$$

где C_{P_i} – мольная теплоемкость i -го вещества.

Целью оптимизации работы химического реактора может являться нахождение максимального выхода целевого продукта R , при этом управляемым параметром будет температура T .

Итак, целевой функцией в этом примере будет функция $C_R(T)$, значения которой при некоторой температуре реактора T определяются в результате численного решения модели (В.2-В.6).

Таким образом, при нелинейной оптимизации очень часто при определении значений целевой функции приходится численно решать задачи, возникающие в сложных моделях. Решению таких задач и будет посвящен наш курс.

Математическая постановка задач оптимизации

Виды ограничений

Несмотря на то, что прикладные оптимизационные задачи относятся к совершенно разным областям, они имеют общую форму. Все эти задачи можно классифицировать как задачи минимизации вещественнозначной функции $f(x)$ на некотором множестве Ω N -мерного векторного аргумента $x = (x_1, x_2, \dots, x_n)$.

Множество Ω задается ограничениями на компоненты вектора x , которые удовлетворяют системе уравнений $h_k(x) = 0$, набору неравенств $g_i(x) \geq 0$, а также ограничены сверху и снизу, т.е. $x_i^{(u)} \geq x_i \geq x_i^{(l)}$.

Иногда множество Ω совпадает со всем N -мерным пространством. В этом случае задача отыскания максимума или минимума, которую также называют задачей оптимизации, называется безусловной.

Заметим, что если функция $f(x)$ имеет в точке x^* минимум, то функция $-f(x)$ в x^* имеем максимум. Поэтому для отыскания максимума применяются те же методы, что и для отыскания минимума. Далее мы, как правило, будем говорить только о задаче отыскания минимума.

Говорят, что функция $f(x)$ имеет локальный минимум x^* , если существует некоторая конечная ε -окрестность точки x^* , в которой выполняется

$$f(x^*) < f(x), \quad |x - x^*| \leq \varepsilon, \quad x \in \Omega. \quad (1.1)$$

У функции может быть много локальных минимумов. Если $f(x^*)$ – наименьший из всех минимумов, то говорят, что функция $f(x)$ достигает абсолютного минимума на множестве Ω . Этот минимум также называют глобальным.

Для нахождения абсолютного минимума надо найти все локальные минимумы, сравнить их и выбрать наименьшее значение.

Поэтому задача отыскания глобального минимума в принципе сводится к задаче (1.1), которую мы и будем рассматривать.

В последующем изложении функцию $f(x)$ будем называть *целевой функцией*,
уравнения $h_k(x) = 0$ – *ограничениями в виде равенств*,
а неравенства $g_i(x) \geq 0$ – *ограничениями в виде неравенств*.

Задача общего вида:

Минимизировать $f(x)$

(пишут $f(x) \rightarrow \min$) при ограничениях

$$h_k(x) = 0, \quad k = 1, \dots, K,$$

$$g_j(x) \geq 0, \quad j = 1, \dots, J,$$

$$x_i^{(u)} \geq x_i \geq x_i^{(l)}, \quad i = 1, \dots, N$$

называется задачей оптимизации с *ограничениями*
или задачей *условной* оптимизации.

Задача, в которой нет ограничений, т.е.

$$J=K=0;$$

$$x_i^{(u)} = -x_i^{(l)} = \infty, \quad i = 1, \dots, N,$$

называется **ОПТИМИЗАЦИОННОЙ** задачей **без ограничений** или задачей **безусловной ОПТИМИЗАЦИИ**.

Критерии оптимальности

Часто оптимизируемая величина связана с экономичностью работы некоторого объекта (аппарат, цех, завод). Оптимизируемый вариант работы объекта должен оцениваться какой-то количественной мерой – *критерием оптимальности*.

Критерием оптимальности называется количественная оценка оптимизируемого качества объекта.

На основании выбранного критерия оптимальности составляется *целевая функция*, представляющая собой зависимость критерия оптимальности от параметров, влияющих на ее значение.

Вид критерия оптимальности или целевой функции определяется конкретной задачей оптимизации.

Таким образом, задача оптимизации сводится к нахождению экстремума целевой функции.

Рассмотрим более подробно требования, которые должны предъявляться к критерию оптимальности.

1. Критерий оптимальности должен выражаться количественно.

2. Критерий оптимальности должен быть единственным.

3. Критерий оптимальности должен отражать наиболее существенные стороны процесса.

4. Желательно чтобы критерий оптимальности имел ясный физический смысл и легко рассчитывался.

Любой оптимизируемый объект схематично можно представить в соответствии с рис. 1.

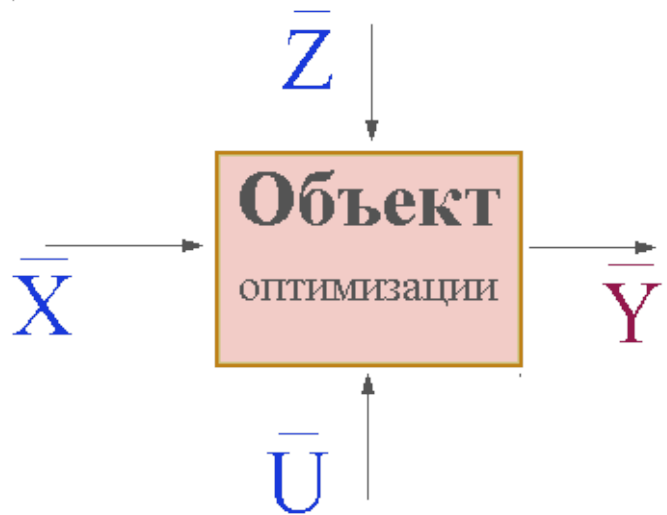


Рис.1. Оптимизируемый объект

\bar{Y} – выходы объекта;
 \bar{X} – контролируемые входные параметры;
 \bar{U} – регулируемые входные параметры (управляемые параметры);
 \bar{Z} – неконтролируемые воздействия;

При постановке конкретных задач оптимизации желательно критерий оптимальности записать в виде аналитического выражения.

В том случае, когда случайные возмущения невелики и их воздействие на объект можно не учитывать, критерий оптимальности может быть представлен как функция *входных*, *выходных* и *управляемых* параметров:

$$R = R(X_1, X_2, \dots, Y_1, Y_2, \dots, U_1, U_2, \dots).$$

Так как $\bar{Y} = F(\bar{U})$, то при фиксированных \bar{X} можно записать:

$$R = R(\bar{U}).$$

При этом всякое изменение значений управляющих параметров двояко сказывается на величине R :

- прямо, так как управляющие параметры непосредственно входят в выражение критерия оптимизации;
- косвенно — через изменение выходных параметров процесса, которые зависят от управляющих.

Если же случайные возмущения достаточно велики и их необходимо учитывать, то следует применять экспериментально-статистические методы, а критерий оптимальности примет вид:

$$R = R(X, U).$$

Процедура решения задачи оптимизации обязательно включает, помимо выбора управляемых параметров, еще и установление ограничений на эти параметры (термостойкость, взрывобезопасность, мощность перекачивающих устройств и т.п.).

Итак, для решения задачи оптимизации необходимо:

а) составить математическую модель объекта оптимизации

$$Y = F(X, U),$$

б) выбрать критерий оптимальности и составить целевую функцию

$$R = \varphi(Y),$$

в) установить возможные ограничения, которые должны накладываться на переменные;

г) выбрать метод оптимизации, который позволит найти экстремальные значения искомых величин.

Классификация задач

Прежде всего, задачи оптимизации можно отнести по типу аргументов к дискретным (компоненты вектора x принимают дискретные или целочисленные значения) и к непрерывным (компоненты вектора x непрерывны).

Для дискретных задач разработаны совершенно специфические методы оптимизации. Мы их рассматривать не будем.

Задачи оптимизации можно классифицировать в соответствии с видом функций f, h_k, g_i и размерностью вектора x .

Задачи, в которых x представляет собой одномерный вектор, называются задачами *с одной переменной* и составляют простейший, но вместе с тем весьма важный подкласс оптимизационных задач.

Задачи условной оптимизации, в которых функции h_k , и g_i являются линейными, носят название задач *с линейными ограничениями*. В таких задачах целевые функции могут быть либо линейными, либо нелинейными.

Задачи, которые содержат только линейные функции вектора непрерывных переменных x , называются *задачами линейного программирования*;

в задачах целочисленного программирования компоненты вектора x должны принимать только целые значения.

Классификацию оптимизационных задач можно представить в следующей форме:

Оптимизация:

Дискретная:

Целочисленное программирование;
Стохастическое программирование.

Непрерывная:

Безусловная:

Глобальная

оптимизация;

Дифференцируемая

оптимизация;

Недифференцируемая

оптимизация.

Условная:

Линейное

программирование;

Нелинейные

задачи;

Стохастическое

программирование.

Более подробная классификация приведена на рис. 2

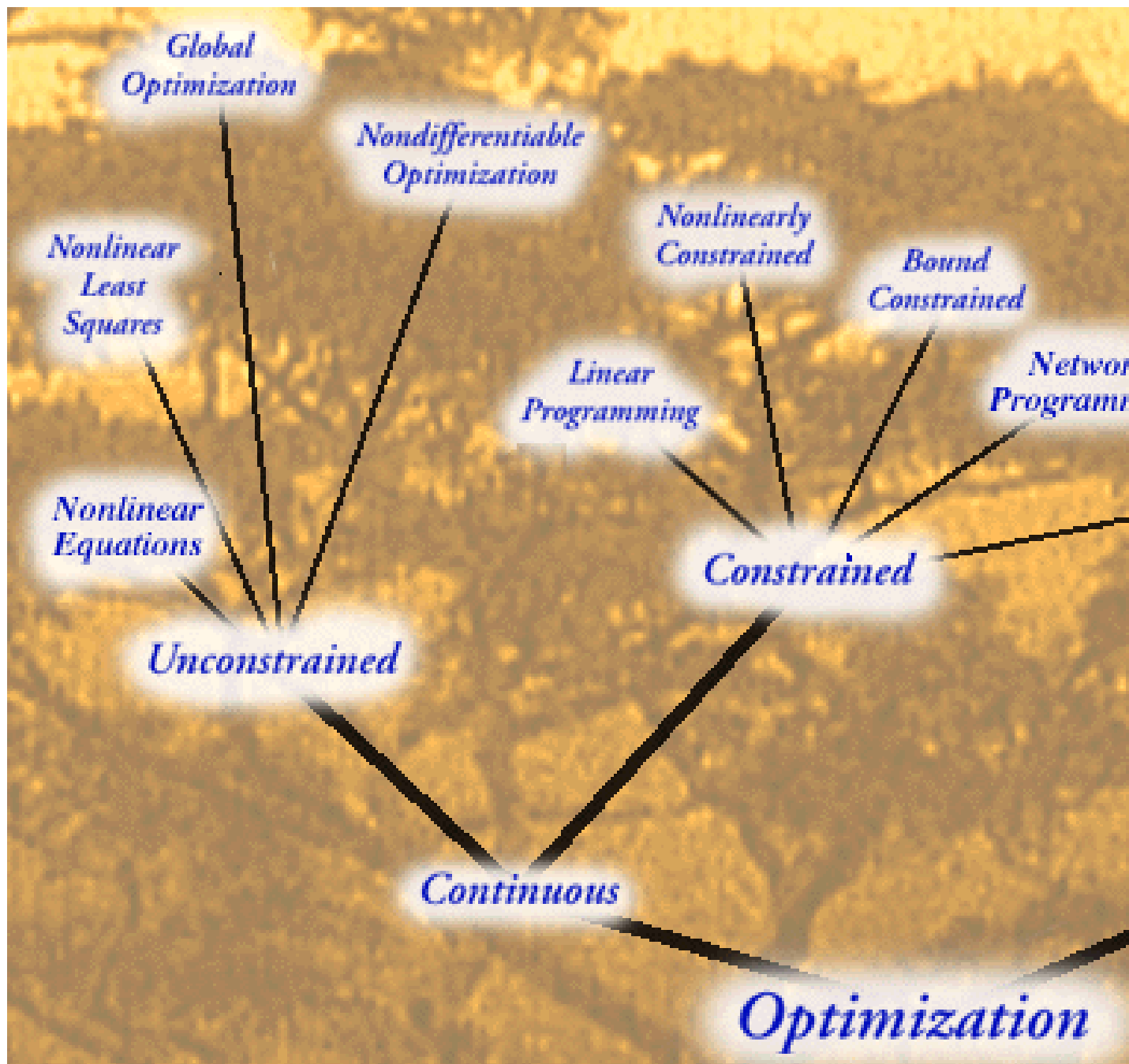
(согласно источнику

<http://www-fp.mcs.anl.gov/otc>

[/_vti_bin/shtml.dll/Guide/](#)

[OptWeb/index.html/map](#)

).



1. Оптимизация в широком смысле слова находит применение в науке, технике и в любой другой области человеческой деятельности.

Оптимизация – целенаправленная деятельность, заключающаяся в получении наилучших результатов при соответствующих условиях.

Поиски оптимальных решений привели к созданию специальных математических методов и уже в 18 веке были заложены математические основы оптимизации (вариационное исчисление, численные методы и др.).

Однако до второй половины 20 века методы оптимизации во многих областях науки и техники применялись очень редко, поскольку практическое использование математических методов оптимизации требовало огромной вычислительной работы, которую без ЭВМ реализовать было крайне трудно, а в ряде случаев – невозможно.

Особенно большие трудности возникали при решении задач оптимизации из-за большого числа параметров и их сложной взаимосвязи между собой.

При наличии ЭВМ ряд задач оптимизации поддается решению.

При постановке задачи оптимизации необходимо:

1. Наличие объекта оптимизации и цели оптимизации. При этом формулировка каждой задачи оптимизации должна требовать экстремального значения лишь одной величины, то есть одновременно системе не должно приписываться два и более критерия оптимизации, так как практически всегда экстремум одного критерия не соответствует экстремуму другого.

Типичный пример неправильной постановки задачи оптимизации:

"Получить максимальную производительность при минимальной себестоимости".

Ошибка заключается в том, что ставится задача поиска оптимума двух величин, противоречащих друг другу по своей сути.

Правильной постановкой задачи может быть:

- а) получить максимальную производительность при заданной себестоимости;
- б) получить минимальную себестоимость при заданной производительности.

В первом случае критерий оптимизации – производительность, а во втором – себестоимость.

2. Наличие ресурсов оптимизации, под которыми понимают возможность выбора значений некоторых параметров оптимизируемого объекта. Объект должен обладать определенными степенями свободы – управляющими воздействиями.

3. Возможность количественной оценки оптимизируемой величины, поскольку только в этом случае можно сравнивать эффекты от выбора тех или иных управляющих воздействий.

4. Учет ограничений.

Пример 1.

Задача о планировании выпуска продукции при ограниченных ресурсах

Нефтеперерабатывающий завод производит за месяц 1500000 л алкилата, 1200000 л крекинг-бензина и 1300000 л изопентола. В результате смешивания этих компонентов в пропорциях 1:1:1 и 3:1:2 получается бензин сорта А и Б соответственно. Стоимость 1000 л бензина сорта А и Б соответственно равна 16000 руб. и 20500 руб.

Определить месячный план производства бензина сорта А и Б, при котором стоимость выпущенной продукции будет максимальной.

Пример 2.

Задача о нахождении размеров нагруженной балки

Задана строительная конструкция, состоящая из балки А длиной $L=35,6$ см и жесткой опоры В. Балка А крепится на жесткой опоре В (рис. 1) с помощью сварного соединения. Балка изготавливается из стали марки 1010 и должна выдержать нагрузку $F = 2721,5$ кг. Размеры h , t , толщину b , ширину l сварного шва необходимо выбрать таким образом, чтобы полные затраты были минимальными.

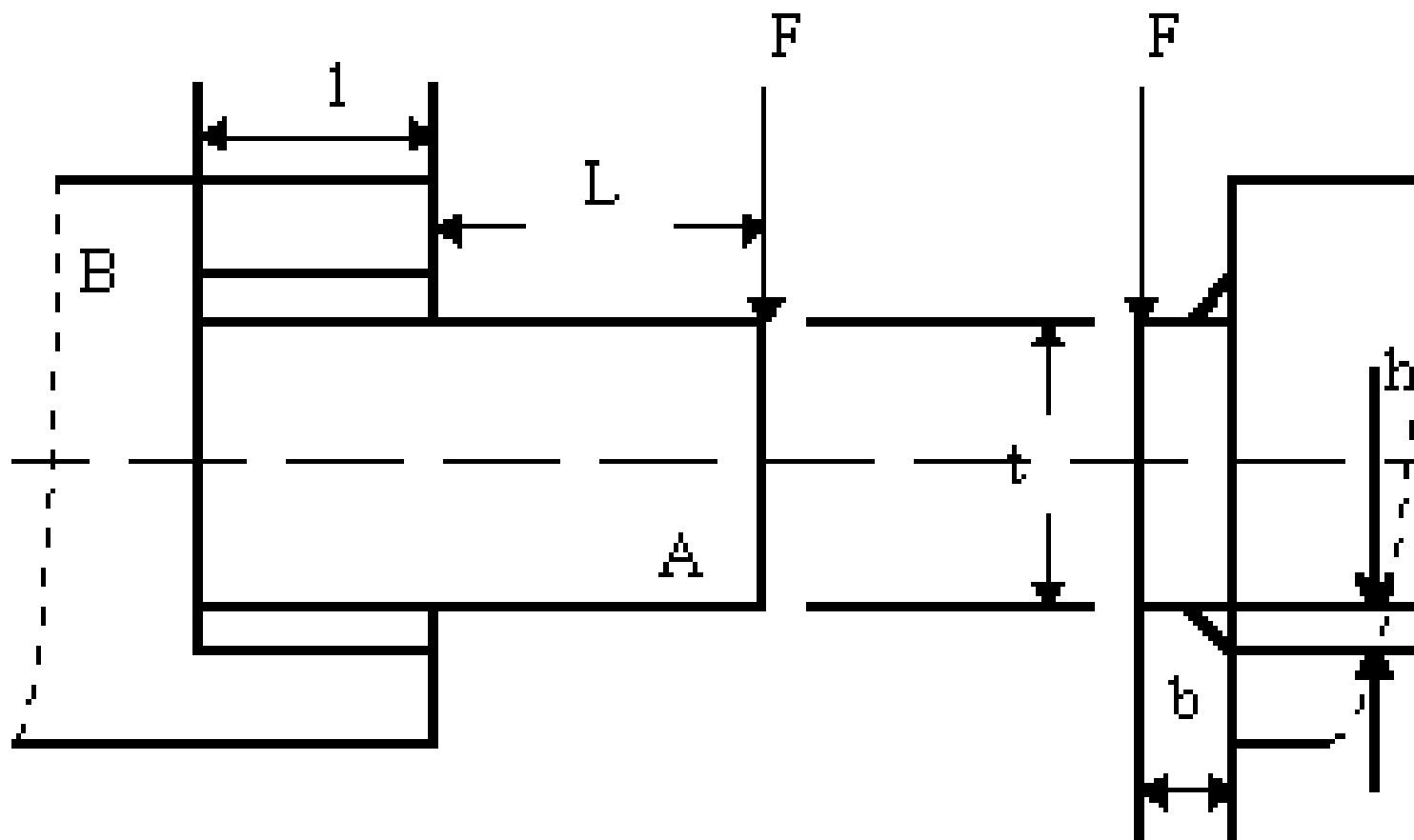


Рис. 1. *Нагруженная балка*

Однако в частных задачах оптимизации, когда объект является частью технологического процесса, не всегда удастся или не всегда целесообразно выделять прямой экономический показатель, который бы полностью характеризовал эффективность работы рассматриваемого объекта. В таких случаях критерием оптимальности может служить технологическая характеристика, косвенно оценивающая экономичность работы агрегата (время контакта, выход продукта, степень превращения, температура). Например, устанавливается оптимальный температурный профиль, длительность цикла "реакция–регенерация".

В задачах оптимизации различают простые и сложные критерии оптимизации.

Критерий оптимальности называется *простым*, если требуется определить экстремум целевой функции без задания условий на какие-либо другие величины. Такие критерии обычно используются при решении частных задач оптимизации (например, определение максимальной концентрации целевого продукта, оптимального времени пребывания реакционной смеси в аппарате и т.п.).

Критерий оптимальности называется **СЛОЖНЫМ**, если необходимо установить экстремум целевой функции при некоторых условиях, которые накладываются на ряд других величин (например, определение максимальной производительности при заданной себестоимости, определение оптимальной температуры при ограничениях по термостойкости катализатора и др.). Процедура решения задачи оптимизации обязательно включает, помимо выбора управляющих параметров, еще и установление ограничений на эти параметры (термостойкость, взрывобезопасность, мощность перекачивающих устройств).

ОДНОМЕРНАЯ ОПТИМИЗАЦИЯ

Оптимизация функции одной переменной – наиболее простой тип оптимизационных задач. Тем не менее, она занимает важное место в теории оптимизации. Это связано с тем, что задачи однопараметрической оптимизации достаточно часто встречаются в инженерной практике и, кроме того, находят свое применение при реализации более сложных итеративных процедур многопараметрической оптимизации.

Разработано большое количество методов одномерной оптимизации и мы рассмотрим две группы таких методов:

- *методы сужения интервала неопределенности;*
- *методы с использованием производных.*

Методы сужения интервала неопределенности

Пусть требуется найти минимум функции $f(x)$ на некотором интервале $[a, b]$. Задача приближенного отыскания минимума в методах сужения интервала неопределенности состоит в том, чтобы найти множество абсцисс x_1, x_2, \dots, x_k , в которых вычисляется функция, такое, что минимальное значение f^* лежит при некотором i в интервале $x_{i-1} \leq x^* \leq x_i$.

Такой интервал называется интервалом неопределенности D . Очевидно, что сначала интервал неопределенности D совпадает с отрезком $[a, b]$.

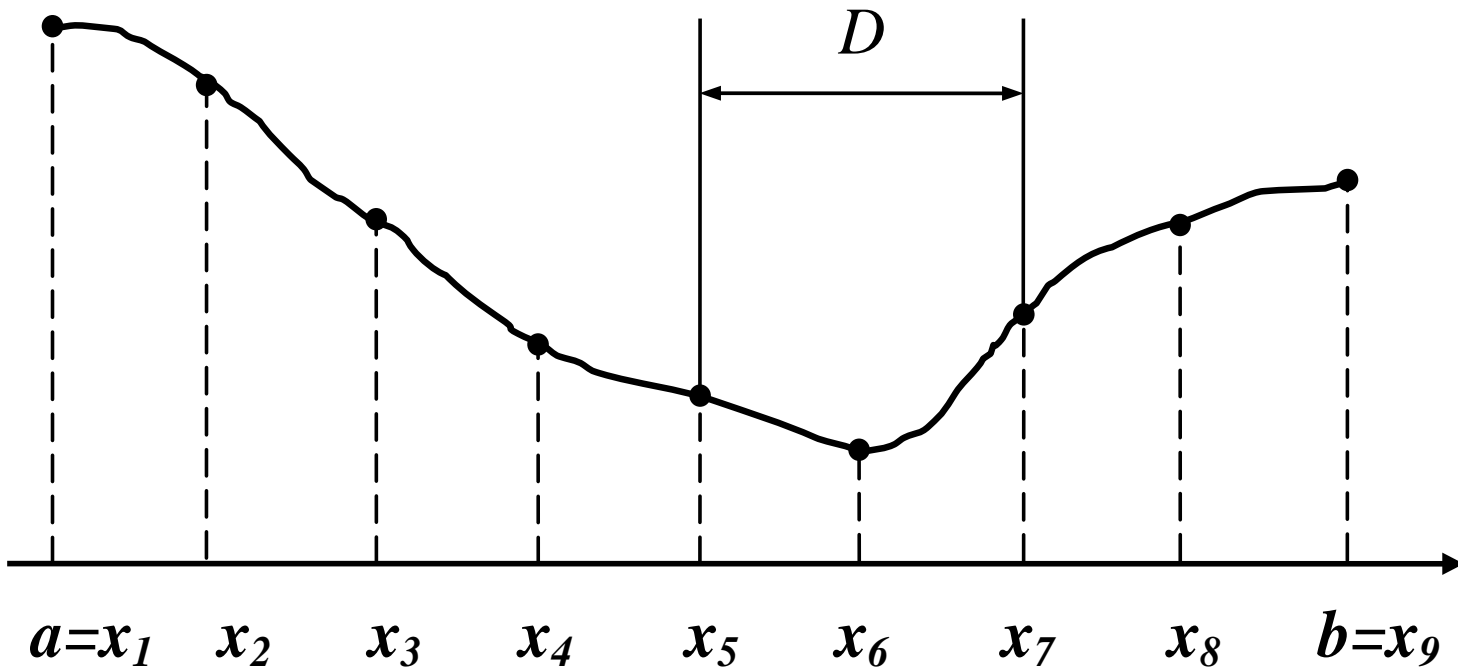
Существуют несколько способов систематического сужения интервала неопределенности.

Рассмотрим три из них.

Общий поиск

Пусть требуется найти минимум функции $f(x)$ на некотором интервале $[a, b]$. Если о функции $f(x)$ на этом интервале никакой дополнительной информации неизвестно, то для поиска минимума на $[a, b]$ можно применить простейший метод перебора, или, иначе, общего поиска.

В этом методе интервал $[a, b]$ делится на несколько равных частей с последующим вычислением значений функции в узлах полученной сетки. В качестве минимума принимается абсцисса с минимальным вычисленным значением функции.



В результате интервал неопределенности сужается до двух шагов сетки.

Обычно говорят о дроблении интервала неопределенности, которое характеризуется коэффициентом α . Разделив интервал неопределенности на n равных частей, получим $n+1$ узел. Тогда $\alpha = \frac{2}{n}$. При этом необходимо вычислить функцию

$N = n+1$ раза.

Следовательно,

$$\alpha = \frac{2}{N-1}, \quad N = 3, 4, 5 \dots \quad (2.1)$$

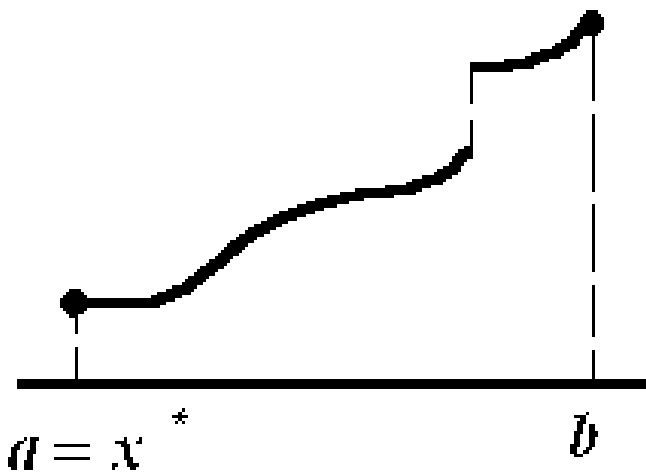
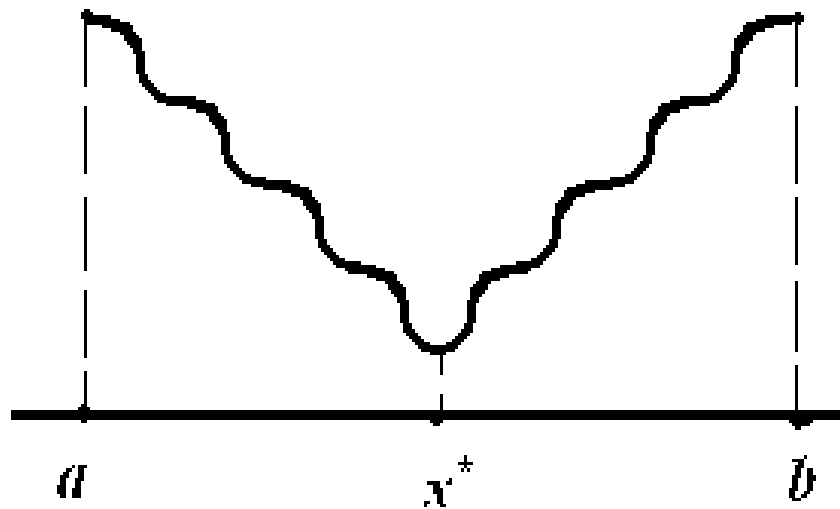
Чтобы получить значение $\alpha = 0.01$ потребуется вычислить функцию в 201 точке, а при $\alpha = 0.001$ $N = 2001$.

Ясно, что эффективность этого метода с уменьшением интервала неопределенности быстро падает.

Унимодальные функции

Более эффективные методы можно построить, если предположить, что исследуемая функция имеет в рассматриваемом интервале только один минимум. Более точно: предположим, что в интервале $[a, b]$ имеется единственное значение x^* такое, что $f(x^*)$ – минимум $f(x)$ на $[a, b]$ и что $f(x)$ строго убывает для $x \leq x^*$ и строго возрастает для $x \geq x^*$. Такая функция называется унимодальной.

Для ее графика имеются три различные формы:

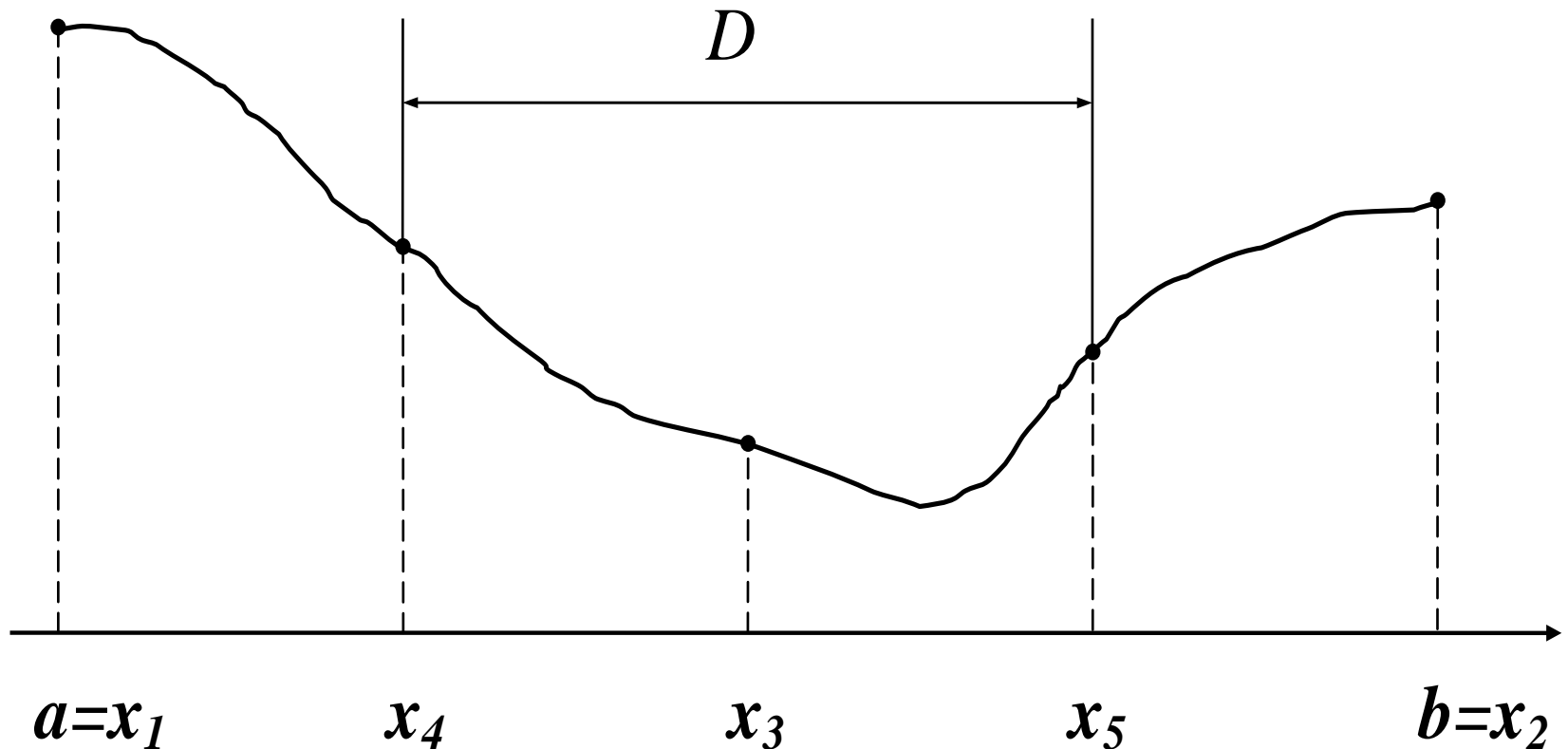


Заметим, что унимодальная функция не обязана быть гладкой или даже непрерывной.

Из предположения унимодальности следует, что для любых точек x_1, x_2 интервала $[a, b]$, таких, что $x_1 < x_2 \leq x^*$ справедливо $f(x_2) < f(x_1)$. Аналогично, если $x^* \leq x_1 < x_2$, то $f(x_2) > f(x_1)$. Обратно, если $x_1 < x_2$ и $f(x_1) > f(x_2)$, то $x_1 \leq x^* \leq b$, а если $f(x_1) < f(x_2)$, то $a \leq x^* \leq x_2$. Далее будем считать исследуемую функцию унимодальной.

Метод деления интервала пополам

Разделим интервал $[a, b]$ на две равные части, а затем каждую из частей еще на две равные части.



Это первый этап поиска минимума. На нем после пяти вычислений функции (два – на краях и три – внутри интервала $[a, b]$) интервал неопределенности сужается вдвое, то есть на этом этапе $\alpha = 0,5$. Новый интервал неопределенности $[x_4, x_5]$ снова разделим пополам, а затем каждую половину снова пополам.

Теперь значения функции по краям и в его середине уже известны.

Поэтому для завершения поиска на этом этапе требуется вычислить только два значения функции, после чего интервал неопределенности снова уменьшится вдвое.

Это преимущество рассмотренного метода сохранится и в дальнейшем.

После N вычислений функции коэффициент дробления интервала составляет

$$\alpha = (0,5)^{\frac{N-3}{2}}, \quad N = 5, 7, 9, \dots \quad (2.2)$$

Здесь $N=5, 7, 9, \dots$, так как интервал неопределенности, начиная со второго этапа, уменьшается только после двух вычислений функции. Из (2.1), (2.2) видно, что метод деления пополам эффективнее, чем общий поиск.

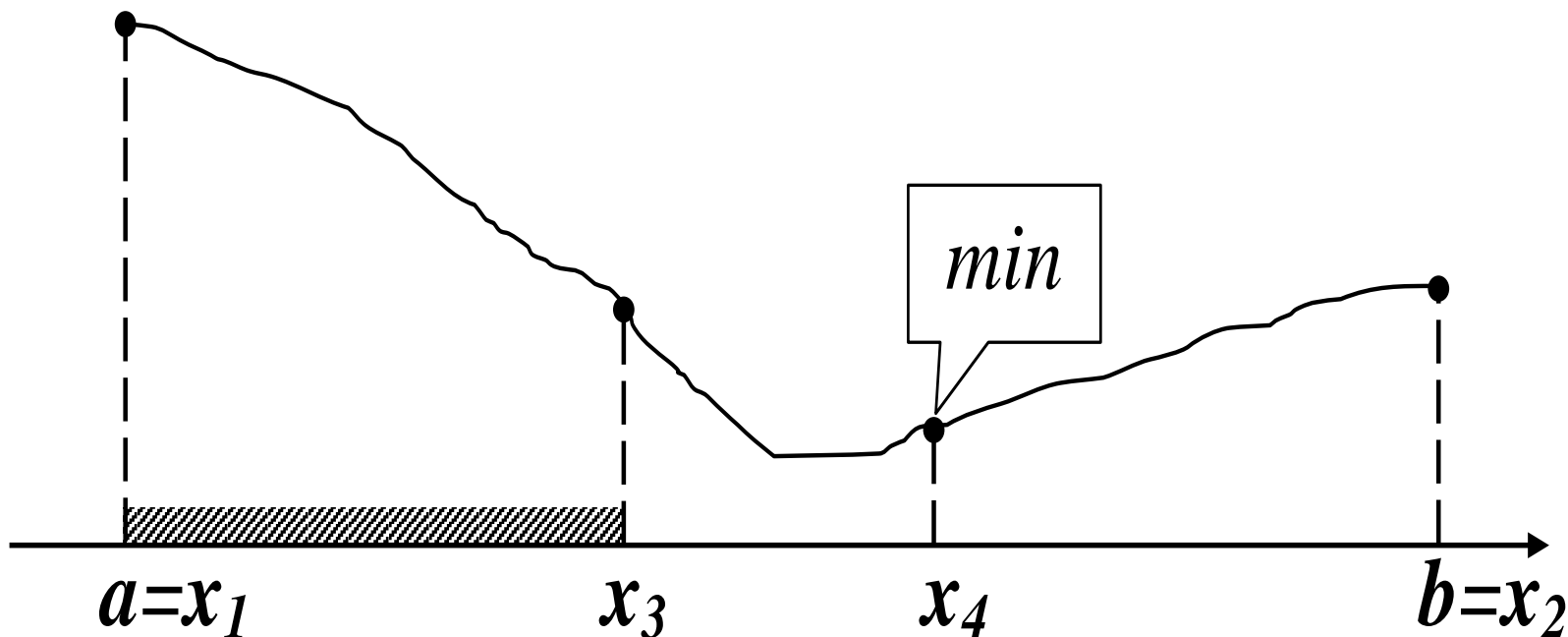
Метод золотого сечения

Деление интервала на неравные части позволяет найти еще более эффективный метод.

Вычислим функцию на концах отрезка $[a, b]$ и положим $a = x_1$, $b = x_2$.

Вычислим также функцию в двух внутренних точках x_3 , x_4 .

Сравним все четыре значения функции и выберем среди них наименьшее.



Пусть, например, наименьшим оказалось $f(x_4)$. Очевидно, минимум находится в одном из прилегающих к нему отрезков.

Поэтому отрезок $[a, x_3]$ можно отбросить и оставить отрезок $[x_3, b]$.

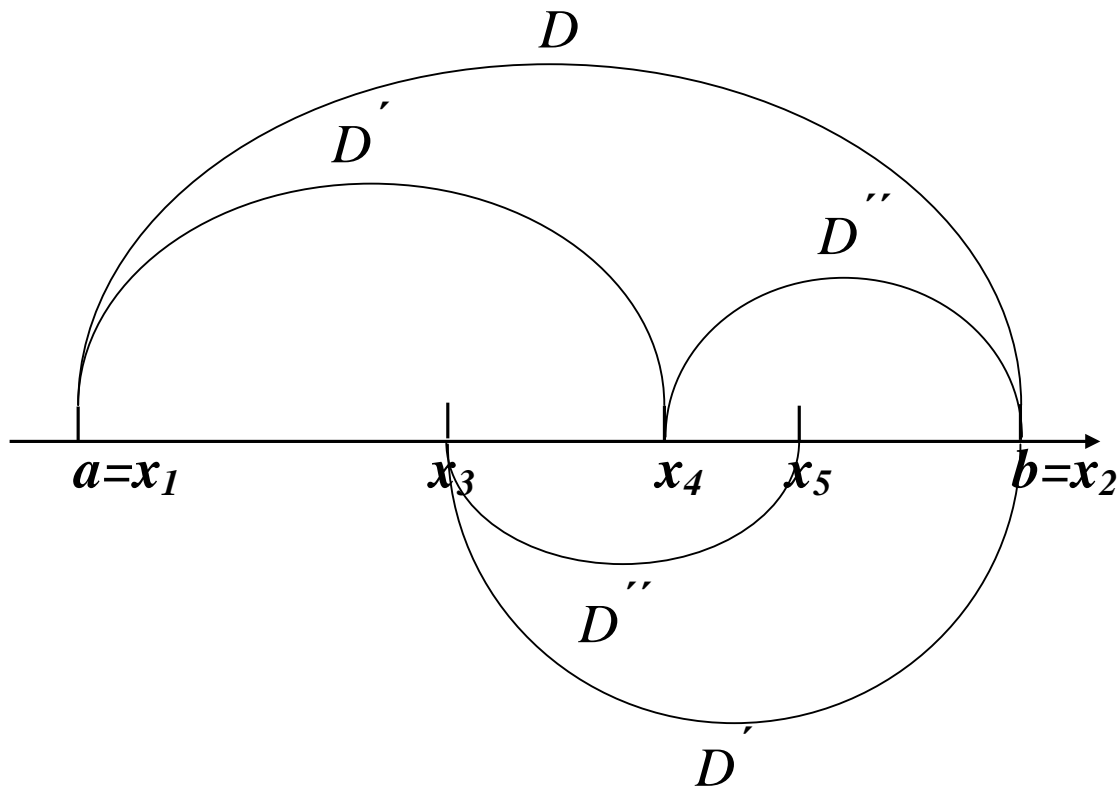
Первый шаг сделан. На отрезке $[x_3, b]$ снова надо выбрать две внутренние точки, вычислив в них и на концах значения функции и сделать следующий шаг.

Но на предыдущем шаге вычислений мы уже нашли функцию на концах нового отрезка $[x_3, b]$ и в одной его внутренней точке x_4 . Потому достаточно выбрать внутри $[x_3, b]$ еще одну точку x_5 , определить в ней значение функции и провести необходимые сравнения. Это вчетверо уменьшает объем вычислений на одном шаге процесса.

Как выгодно размещать точки?

Каждый раз оставшийся отрезок делится на три части и затем отбрасывается один из крайних отрезков.

Обозначим первоначальный интервал неопределенности через D .



Так как в общем случае может быть отброшен любой из отрезков x_1x_3 или x_4x_2 , то выберем точки x_3 и x_4 так, чтобы длины этих отрезков были одинаковы:

$$x_3 - x_1 = x_2 - x_4.$$

После отбрасывания получится новый интервал неопределенности длины D' .

Обозначим отношение $\frac{D}{D'}$ буквой φ :

$$\varphi = \frac{D}{D'}.$$

Далее продолжим процесс аналогично.

Для этого интервал D' разделим подобно

интервалу D , то есть положим $\frac{D'}{D''} = \frac{D}{D'} = \varphi$, где

D'' – следующий интервал неопределенности.

Но D'' по длине равен отрезку, отброшенному на предыдущем этапе, то есть $D'' = D - D'$.

Поэтому получим:

$$\frac{D}{D'} = \frac{D'}{D - D'} \Rightarrow \frac{D'}{D} = \frac{D}{D'} - 1.$$

Это приводит к уравнению $\frac{1}{\varphi} = \varphi - 1$ или,

что то же

$$\varphi^2 - \varphi - 1 = 0.$$

Положительный корень этого уравнения дает

$$\varphi = \frac{\sqrt{5} + 1}{2} \approx 1,6180.$$

Последнее число известно в математике как золотое отношение, а описанное деление отрезка как золотое сечение.

Потому рассматриваемый метод поиска минимума называют методом золотого сечения.

Отношение $\frac{D}{D'} = \varphi \approx 1,618$ показывает, во сколько раз сокращается интервал неопределенности при одном добавочном вычислении функции.

Учтем, что первые три вычисления еще не сокращают интервал неопределенности.

Поэтому после N вычислений функции коэффициент дробления будет

$$\alpha = \left(\frac{1}{\varphi} \right)^{N-3} \approx (0,6180)^{N-3}. \quad (2.3)$$

При $N \rightarrow \infty$ длина интервала неопределенности стремится к нулю как геометрическая прогрессия со знаменателем $\frac{1}{\varphi}$, то есть метод золотого сечения всегда сходится.

Очевидно, этот метод более эффективен, чем метод деления пополам, так как после N вычислений функции длина интервала неопределенности уменьшается при золотом сечении в $\varphi^{N-3} \approx (1,6180)^{N-3}$ раз, а в методе деления пополам в $2^{\frac{n-3}{2}} \approx (1,4142)^{N-3}$ раза.

Приведем теперь вычислительную схему метода.

Имеем $\frac{D}{D'} = \varphi$, причем

$$D = x_2 - x_1, \quad D' = x_4 - x_1 \quad \text{или} \quad x_2 - x_3.$$

Поэтому $\frac{x_2 - x_1}{x_2 - x_3} = \varphi, \quad \frac{x_2 - x_1}{x_4 - x_1} = \varphi,$

что дает:

$$x_3 = x_2 - \frac{1}{\varphi}(x_2 - x_1) = x_2 - \frac{\sqrt{5}-1}{2}(x_2 - x_1) \approx x_2 - 0,6180(x_2 - x_1) \quad (2.4)$$

$$x_4 = x_1 + \frac{1}{\varphi}(x_2 - x_1) = x_1 + \frac{\sqrt{5}-1}{2}(x_2 - x_1) \approx x_1 + 0,6180(x_2 - x_1) \quad (2.5)$$

Так как длины отрезков x_1x_3 и x_4x_2 равны, последнее равенство можно переписать следующим образом:

$$x_4 = x_1 + x_2 - x_3. \quad (2.6)$$

После сравнения может быть отброшена точка с любым номером, так что на следующих шагах оставшиеся точки будут перенумерованы беспорядочно.

Пусть на данном отрезке есть четыре точки, x_i, x_j, x_k, x_l , из которых какие-то две являются концами отрезка.

Выберем ту точку, в которой функция принимает наименьшее значение; пусть это оказалась точка x_i :

$$f(x_i) < f(x_j), f(x_k), f(x_l) \quad (2.7)$$

Затем отбрасываем ту точку, которая более удалена от x_i (это верно в методе золотого сечения). Пусть этой точкой оказалась x_l :

$$|x_l - x_i| > |x_j - x_i|, |x_k - x_i|.$$

Определим порядок распределения оставшихся трех точек на числовой оси; пусть, например, $x_k < x_i < x_j$. Пронумеруем эти точки, положив $k=1$, $j=2$, $i=3$.

Тогда новую внутреннюю точку введем по формуле (2.6):

$$x_4 = x_1 + x_2 - x_3.$$

Ее номер теперь – 4.

Вычислим функцию $f(x_4)$ в этой точке. Выполним сравнение, отбросим одну точку, заново переименуем точки, введем новую точку по формуле (2.6) и т.д.

Минимум находится где-то внутри последнего отрезка: $x^* \in [x_1, x_2]$.

Поэтому процесс прекращается, когда длина этого интервала неопределенности станет меньше заданной погрешности:
 $x_2 - x_1 < \varepsilon$.

Заметим, что если на $[a, b]$ функция имеет несколько минимумов, то процесс сойдется к одному из них, но не обязательно к наименьшему.

Приведем таблицу сравнения методов поиска минимума по значениям коэффициента дробления интервала неопределенности после N вычислений функции:

N	Коэффициент дробления α		
	Общий поиск	Деление пополам	Золотое сечение
3	1	1	1
4	0,667	-	0,618
5	0,500	0,500	0,382
6	0,400	-	0,250
7	0,333	0,250	0,146
8	0,286	-	0,090
9	0,250	0,125	0,056
10	0,222	-	0,0345
19	0,111	0,00391	0,000453
20	0,105	-	0,000280
21	0,100	0,00195	0,000173

Установление первоначального интервала неопределенности

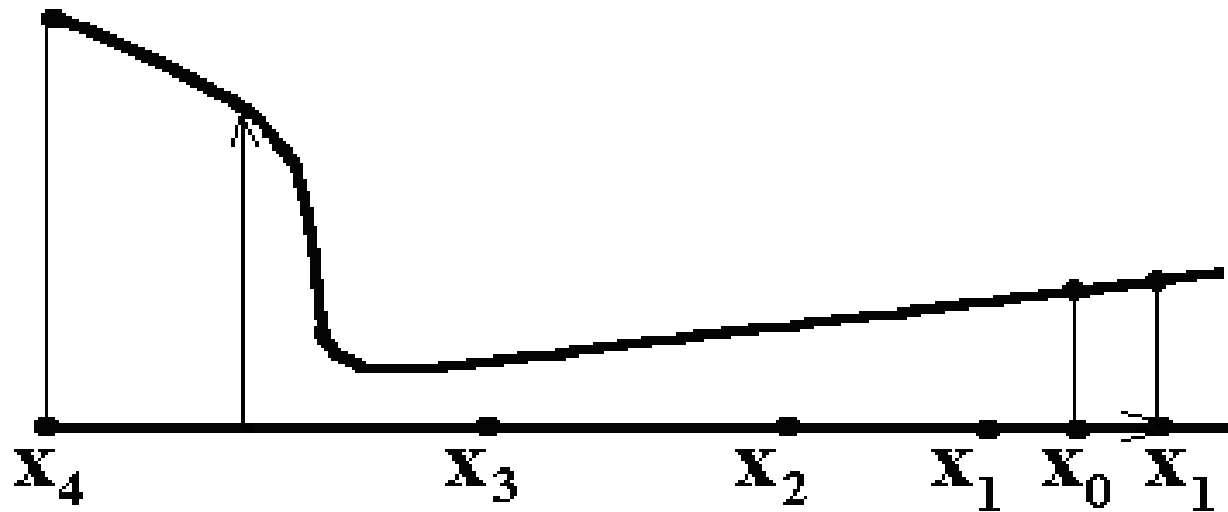
Рассмотренные выше методы поиска минимума, которые позволяют определить оптимум функции одной переменной путем уменьшения интервала поиска, носят название методов исключения интервалов.

Процесс применения методов поиска на основе исключения интервалов включает два этапа:

- этап установления границ интервала;
- этап уменьшения интервала.

Рассмотрим теперь этап установления границ интервала. Обычно используется эвристический метод, например, метод Свенна.

Итак, пусть требуется найти минимум функции $f(x)$ не на отрезке, а на всей оси x . Предположим снова, что функция $f(x)$ унимодальна. Выберем некоторое начальное приближение x_0 , и сделаем из него шаг некоторой длины h : $x_1 = x_0 + h$. Если $f(x_1)$ окажется большее, чем $f(x_0)$, то изменим направление шага и положим $x_1 = x_0 - h$. Пусть теперь $f(x_1) < f(x_0)$. Удвоим шаг $h' = 2h$ и положим $x_2 = x_1 + h'$ и т.д., до тех пор, пока на некотором шаге не будет выполнено условие $f(x_n) > f(x_{n-1})$.



Теперь ясно, что минимум унимодальной функции лежит на отрезке $[x_4, x_3]$ и его можно найти одним из рассмотренных методов.

Рассмотренные методы оптимизации используют только значения функции . Такие методы называются методами 0-го порядка. Скорость их сходимости невелика. Если предположить, что функция дифференцируема, то можно предложить более быстрые методы, использующие производные. Методы, использующие *первую* производную, называются методами 1-го порядка и т. д.

Ньютоновские методы

Пусть функция $f(x)$ дважды дифференцируема. Как известно из математического анализа, условием минимума такой функции является равенство

$$f'(x^*) = 0. \quad (2.8)$$

Это необходимое условие. Для того чтобы точка x^* была минимумом, достаточно выполнения условия

$$f''(x^*) > 0. \quad (2.9)$$

Будем численно решать уравнение

$$f'(x) = 0. \quad (2.10)$$

Зададим некоторое начальное приближение x_k , разложим в этой точке функцию в ряд Тейлора и ограничимся лишь членами до второго порядка включительно, т.е. построим квадратичную модель функции:

$$\hat{f}(x) \approx f(x_k) + f'(x_k)(x - x_k) + \frac{1}{2} f''(x_k)(x - x_k)^2. \quad (2.11)$$

Если $f''(x_k) \neq 0$, $\hat{f}(x)$ будет иметь единственную стационарную точку. Найдем ее, для чего приравняем нулю производную $\hat{f}'(x)$:

$$\hat{f}'(x) = f'(x_k) + f''(x_k)(x - x_k) = 0.$$

Решим это уравнение относительно x и найденное решение примем за очередное, $k+1$ -ое приближение к минимуму:

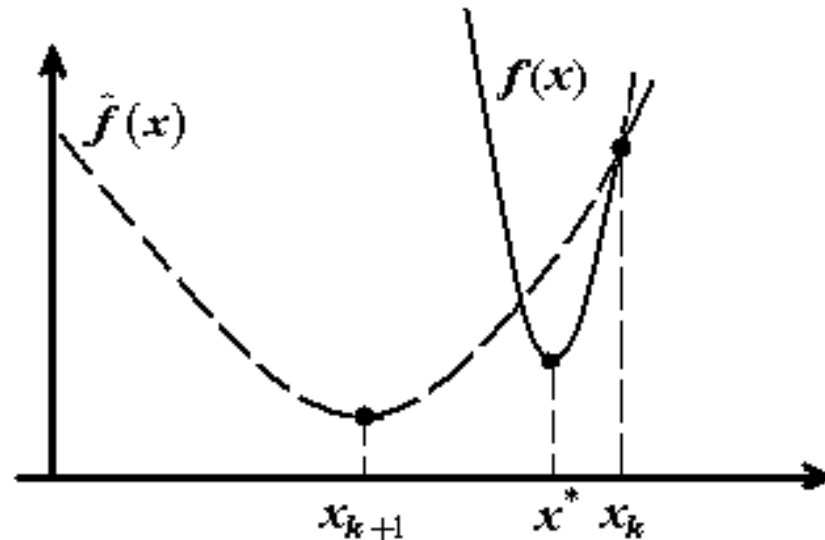
$$x_{k+1} = x_k - \frac{f'(x_k)}{f''(x_k)}. \quad (2.12)$$

Формулу (2.12) можно получить иначе, если применить численный метод решения уравнения $g(x) = 0$, известный, как метод Ньютона:

$$x_{k+1} = x_k - \frac{g(x_k)}{g'(x_k)}.$$

Надо только учесть, что теперь мы решаем уравнение $f'(x) = 0$, то есть положить $g(x) = f'(x)$.

У алгоритма (2.12) есть два недостатка. Во-первых, уравнение $f'(x) = 0$ может определять не только минимум, но и максимум. Во-вторых, модельная функция $\hat{f}(x)$ может сильно отличаться от оптимизируемой функции $f(x)$ и шаг $x_{k+1} - x_k$ может оказаться слишком большим:



Поэтому стратегию (2.12) следует уточнить.

Чтобы быть уверенными, что мы продвигаемся к минимуму, будем на каждом шаге проверять соотношение $f(x_{k+1}) < f(x_k)$.

Если оно выполняется, то переходим к следующему шагу и т.д. Если же $f(x_{k+1}) > f(x_k)$, а $f'(x_k)(x - x_k) < 0$, то функция $f(x)$ должна первоначально уменьшаться в направлении от x_k к x_{k+1} , поэтому следующую приемлемую точку можно найти, дробя шаг в обратном направлении, например, ПОЛОЖИВ

$$x'_{k+1} = \frac{x_{k+1} + x_k}{2}.$$

Из формулы (2.12) видно, что выражение $f'(x_k)(x - x_k)$ отрицательно тогда и только тогда, когда $f''(x_k)$ положительна. Это означает, что если локальная модель, используемая для получения Ньютоновского шага, имеет минимум, а не максимум, то гарантируется существование подходящего направления шага.

С другой стороны, если $f''(x_k) < 0$ и $f'(x_k)(x - x_k) > 0$, то при переходе от x_k к x_{k+1} , $f(x)$ первоначально увеличивается, поэтому шаг нужно сделать в противоположном направлении.

Критерий останова для оптимизации можно выбрать в виде

$$\left| \frac{f'(x_{k+1})}{f(x_{k+1})} \right| < \varepsilon, \quad (2.13)$$

где ε – заранее заданная точность.

Описанный метод с основным шагом (2.12) и приведенными уточнениями обычно называют методом Ньютона или *Ньютона-Рафсона*.

В некоторых задачах производные функции $f(x)$ недоступны и метод Ньютона можно модифицировать.

Выберем начальное приближение x_k и малый шаг h .

Рассмотрим три точки $x_k - h$, x_k , $x_k + h$. Тогда производные $f'(x_k)$ и $f''(x_k)$ можно аппроксимировать следующим образом:

$$f'(x_k) = \frac{f(x_k + h) - f(x_k - h)}{2h},$$

$$f''(x_k) = \frac{f'(x_k + h) - f'(x_k - h)}{2h} =$$

$$= \frac{\frac{f(x_k + h) - f(x_k)}{h} - \frac{f(x_k) - f(x_k - h)}{h}}{2h} =$$

$$= \frac{f(x_k + h) - 2f(x_k) + f(x_k - h)}{2h^2}.$$

Подставляя это в алгоритм (2.12), найдем:

$$x_{k+1} = x_k - h \frac{f(x_k + h) - f(x_k - h)}{f(x_k + h) - 2f(x_k) + f(x_k - h)}. \quad (2.14)$$

Формула (2.14) дает основной шаг алгоритма, называемого квазиньютоновским методом или модифицированным методом Ньютона.

Все соображения относительно шага $x_{k+1} - x_k$, приводимые при выводе метода ***Ньютона-Рафсона***, остаются в силе.

МИНИМУМ ФУНКЦИИ МНОГИХ ПЕРЕМЕННЫХ

Рельеф функции

Понятие "рельеф функции" удобно рассмотреть на примере функции двух переменных $z = F(x, y)$. Эта функция описывает некоторую поверхность в трехмерном пространстве с координатами x, y, z . Задача $F(x, y) \rightarrow \min$ означает поиск низшей точки этой поверхности.

Как в топографии, изобразим рельеф этой поверхности *линиями уровня*.

Проведем равноотстоящие плоскости
 $z = \text{const}$

и найдем линии их пересечения с поверхностью
 $F(x, y)$.

Проекции этих линий на плоскость x, y
называют *линиями уровня*.

Направление убывания функции будем указывать штрихами рядом с линиями уровня. По виду линий уровня условно выделим три типа рельефа:

котловинный,
овражный,
неупорядоченный.

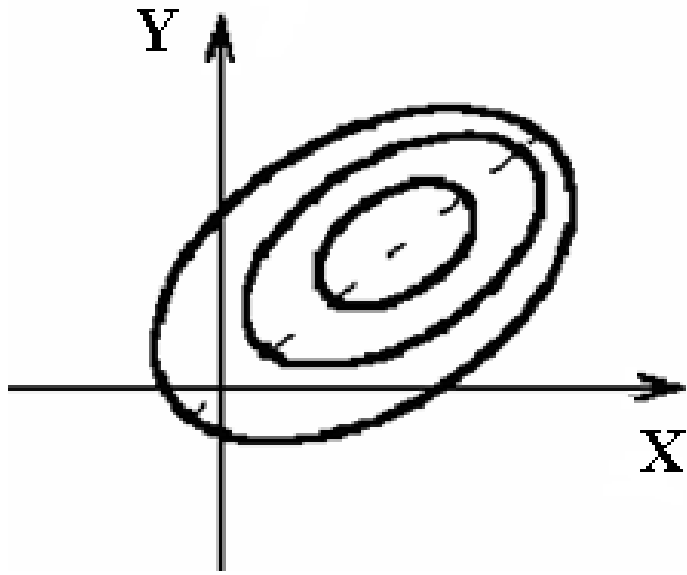


Рис.4. Котловинный рельеф
рельеф

При котловинном рельефе линии уровня похожи на эллипсы (Рис. 4).

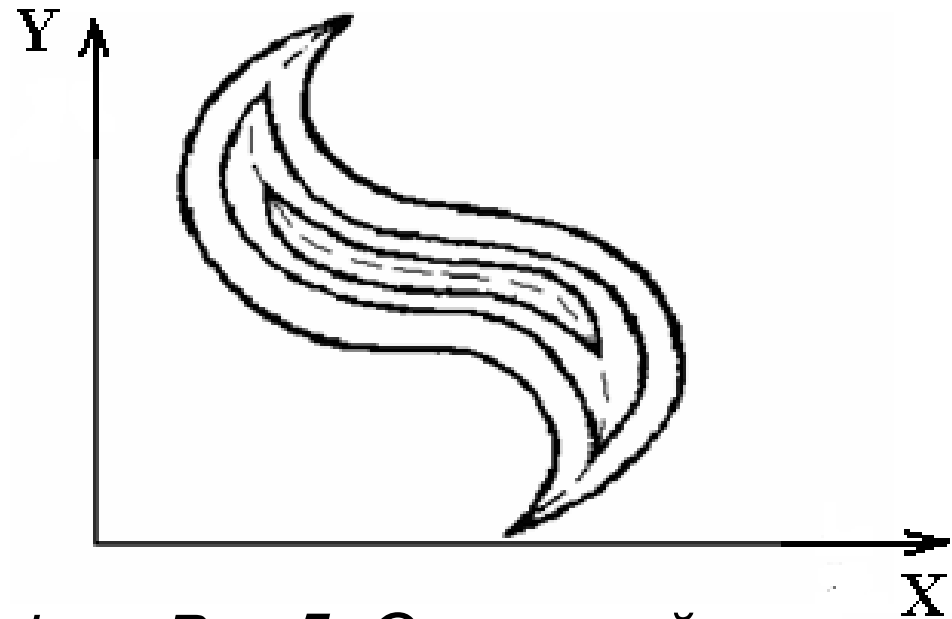


Рис.5. Овражный

Рассмотрим овражный тип рельефа.
Если линии уровня кусочно-гладкие (Рис. 5),
то выделим на каждой из них точку излома.
Геометрическое место точек излома назовем
истинным оврагом, если угол направлен в
сторону возрастания функции,
и *гребнем*, – если в сторону убывания.

Чаще линии уровня всюду гладкие, но на них имеются участки с большой кривизной. Геометрические места точек с наибольшей кривизной назовем *разрешимым оврагом* или *гребнем* (Рис. 6).

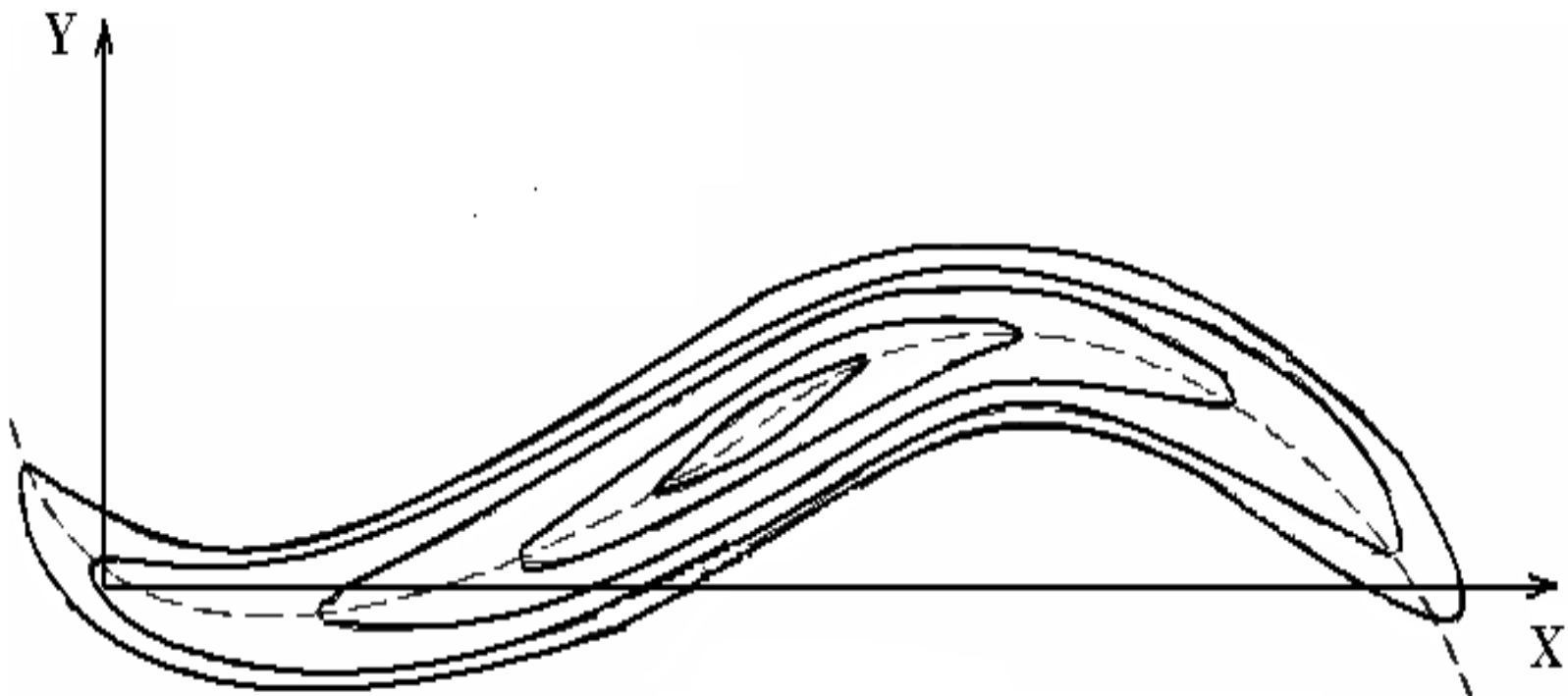


Рис. 6. Разрешимый овраг

Например, рельеф функции

$$F(x, y) = 10(y - \sin x)^2 + 0.1x^2$$

(Рис. 7) имеет ярко выраженный извилистый разрешимый овраг, "дно" которого – синусоида, а низшая точка – начало координат.

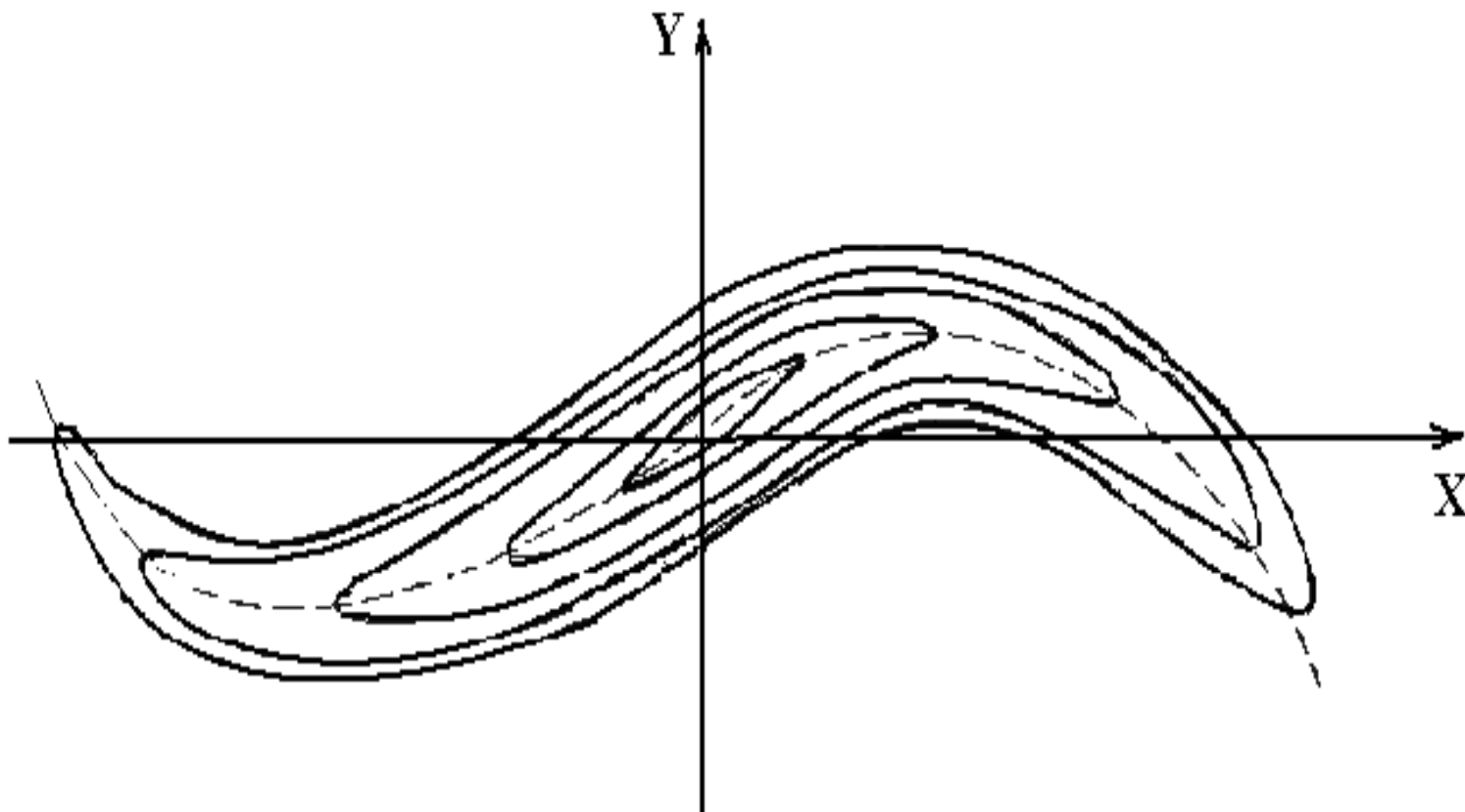


Рис. 7. Рельеф функции

$$F(x, y) = 10(y - \sin x)^2 + 0.1x^2$$

Неупорядоченный тип рельефа характеризуется наличием многих максимумов и минимумов.

Так, функция

$$F(x, y) = (1 + \sin^2 x)(1 + \sin^2 y)$$

(Рис. 8) имеет минимумы в точках $x_k^* = \pi k$, $y_l^* = \pi l$ и максимумы в точках, сдвинутых относительно минимумов на $\frac{\pi}{2}$ по каждой координате.

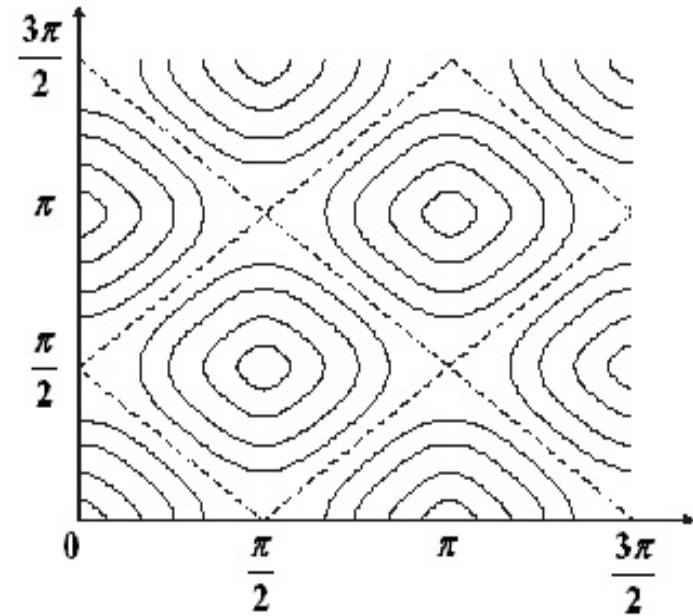
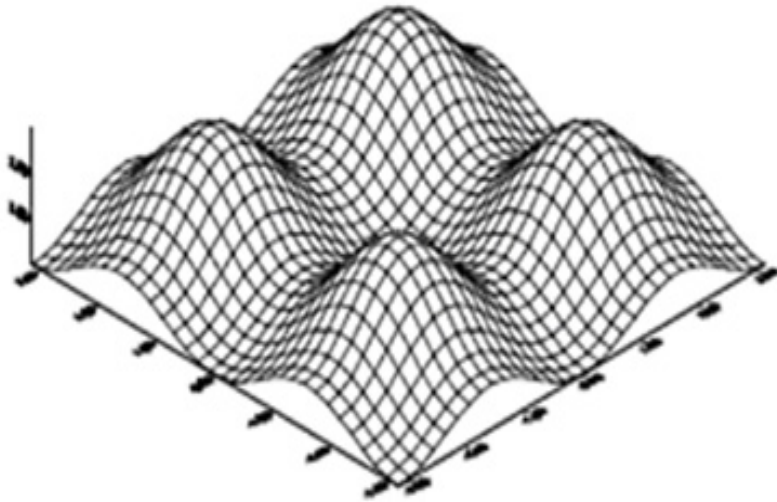


Рис. 8. Функция

$$F(x, y) = (1 + \sin^2 x)(1 + \sin^2 y)$$

и её рельеф

Все эффективные методы поиска минимума сводятся к построению траекторий, вдоль которых функция убывает; разные методы отличаются способами построения таких траекторий.

Метод, приспособленный к одному типу рельефа, может оказаться плохим на рельефе другого типа

Метод покоординатного спуска (Метод Гаусса)

Изложим этот метод на примере функции трех переменных $F(x, y, z)$.

Выберем нулевое приближение x_0, y_0, z_0 .

Фиксируем значение двух координат

$$y = y_0, \quad z = z_0.$$

Тогда функция будет зависеть только от одной переменной x ; обозначим ее через

$$f_1(x) = F(x, y_0, z_0).$$

Используя какой-либо способ одномерной оптимизации, отыщем минимум функции $f_1(x)$ и обозначим его через x_1 .

Мы сделали шаг из точки (x_0, y_0, z_0) в точку (x_1, y_0, z_0) по направлению, параллельному оси x ; на этом шаге значение функции уменьшилось.

Теперь из новой точки сделаем спуск по направлению, параллельному оси y , то есть рассмотрим функцию $f_2(y) = F(x_1, y, z_0)$, найдем ее минимум и обозначим его через y_1 . Вторым шагом приводит нас в точку (x_1, y_1, z_0) .

Из этой точки делаем третий шаг – спуск параллельно оси z и находим минимум функции $f_3(z) = F(x_1, y_1, z)$.

Приход в точку (x_1, y_1, z_1) завершает цикл спусков или первую итерацию.

Далее будем повторять циклы.

На каждом спуске функция не возрастает, и при этом значение функции ограничено снизу ее значением в минимуме $F^* = F(x^*, y^*, z^*)$. Следовательно, итерации сходятся к некоторому пределу $\tilde{F} \geq F^*$.

Будет ли здесь иметь место равенство, то есть сойдутся ли спуски к минимуму и как быстро?

Это зависит от функции и выбора нулевого приближения.

На примере функции двух переменных легко убедиться, что существуют случаи сходимости спуска по координатам к минимуму и случаи, когда этот спуск к минимуму не сходится.

В самом деле, рассмотрим геометрическую трактовку этого метода на рисунках 9 и 10, приведенных ниже:

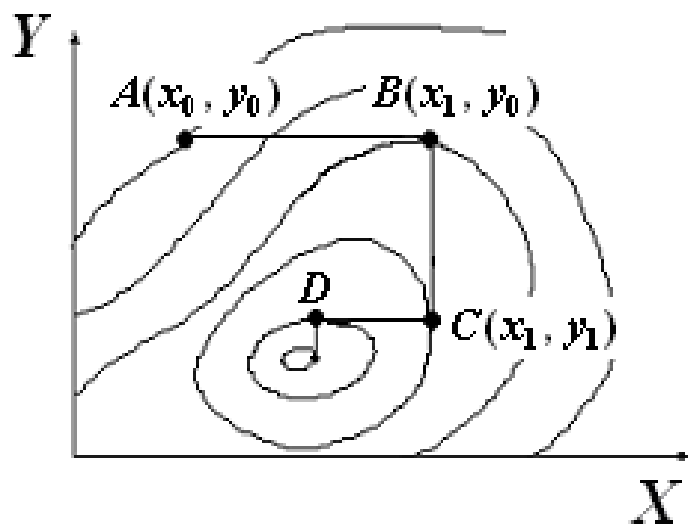


Рис. 9

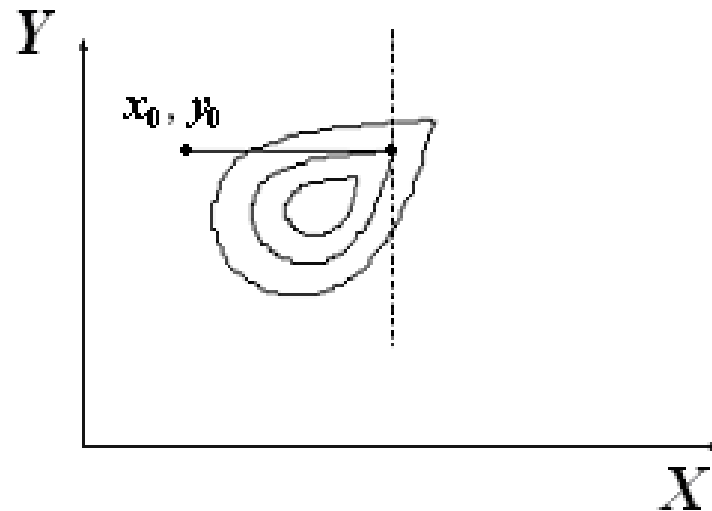


Рис. 10

Будем двигаться по выбранному направлению, то есть по некоторой прямой в плоскости .

В тех участках, где прямая пересекает линии уровня, мы при движении переходим от одной линии уровня к другой, так что при этом движении функция меняется (возрастает или убывает).

Только в той точке, где данная прямая касается линии уровня (рис. 9), функция имеет экстремум вдоль этого направления. Найдя такую точку, мы завершаем в ней спуск по первому направлению и должны начать спуск по второму.

Пусть линии уровня образуют истинный овраг. Тогда возможен случай (рис. 10), когда спуск по одной координате приводит нас на "дно оврага", а любое движение по следующей координате (пунктирная линия) ведет нас на подъем:

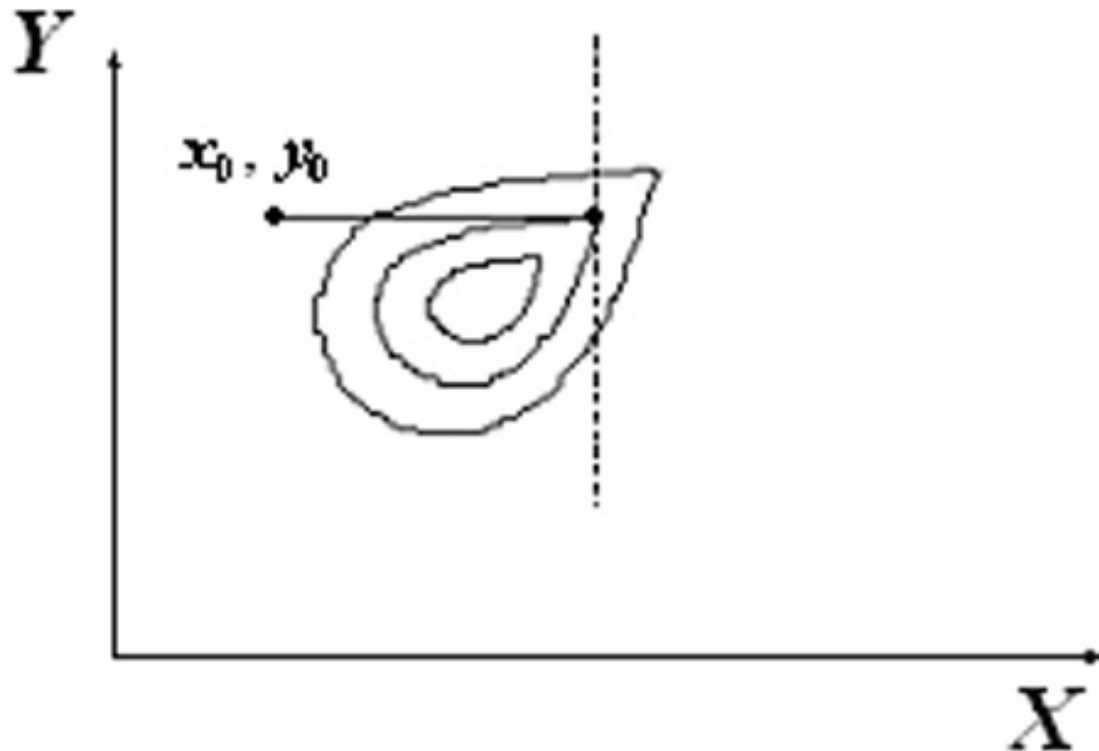


Рис. 10

Никакой дальнейший спуск по координатам невозможен, хотя минимум еще не достигнут.

В данном случае процесс спуска по координатам не сходится к минимуму.

Наоборот, если функция достаточно гладкая, то в некоторой окрестности минимума процесс спуска по координатам сходится к этому минимуму.

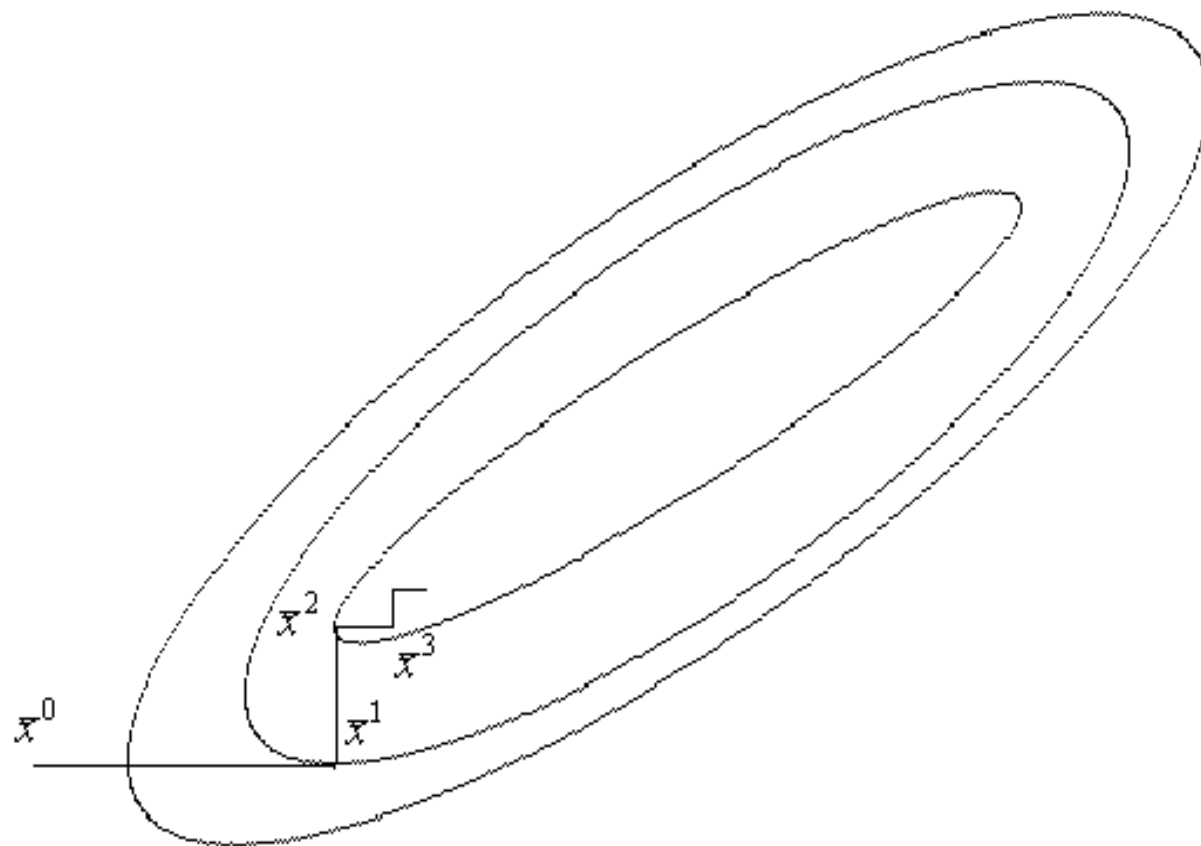


Рис. 9а.

Однако скорость сходимости сильно зависит от формы линий уровня.

Так, если рельеф функции имеет тип "разрешимый овраг", то при попадании траектории спуска в такой овраг сходимость становится настолько медленной, что расчет практически вести невозможно.

Иногда метод покоординатного спуска используют в качестве первой попытки при нахождении минимума.

Метод оврагов

Выберем произвольную точку ρ_0 и спустимся из нее (например, по координатам), делая не очень много шагов, то есть, не требуя высокой точности.

Конечную точку спуска обозначим r_0 .

Если рельеф овражный, эта точка окажется вблизи дна оврага (рис. 11а):

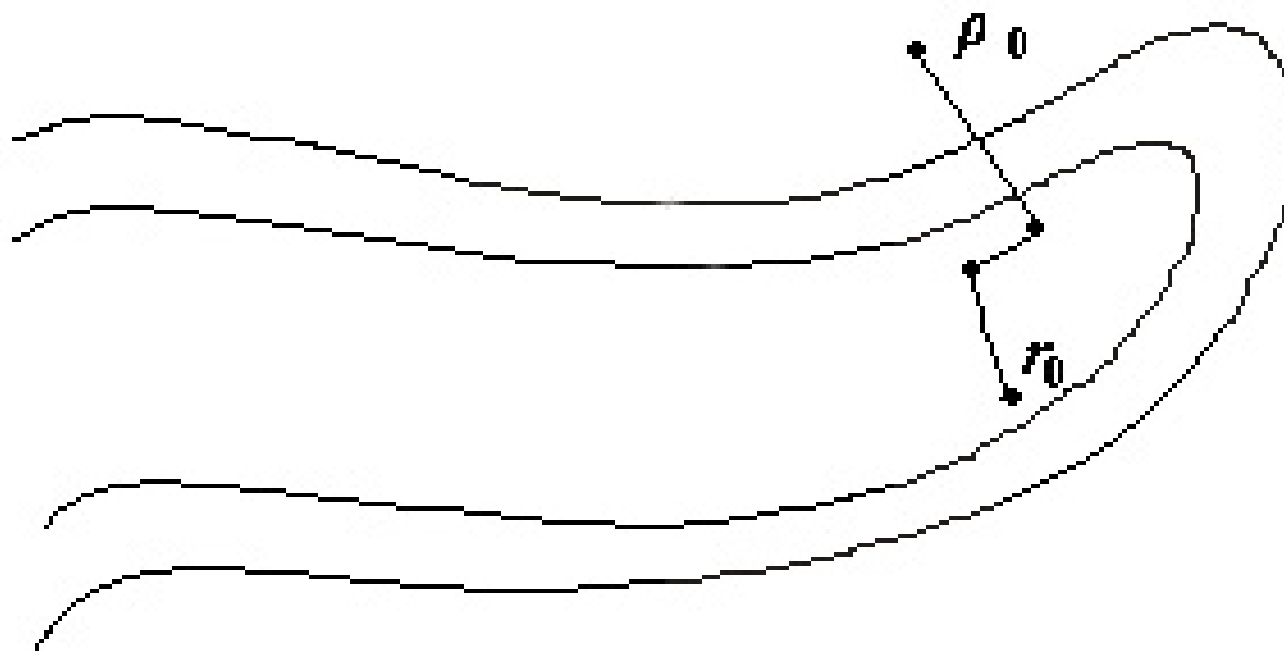


Рис. 11а. Иллюстрация к методу оврагов

Теперь выберем другую точку ρ_1
не слишком далеко от первой.
Из нее также сделаем спуск и
попадем в некоторую точку r_1 .
Эта точка тоже лежит вблизи дна
оврага.

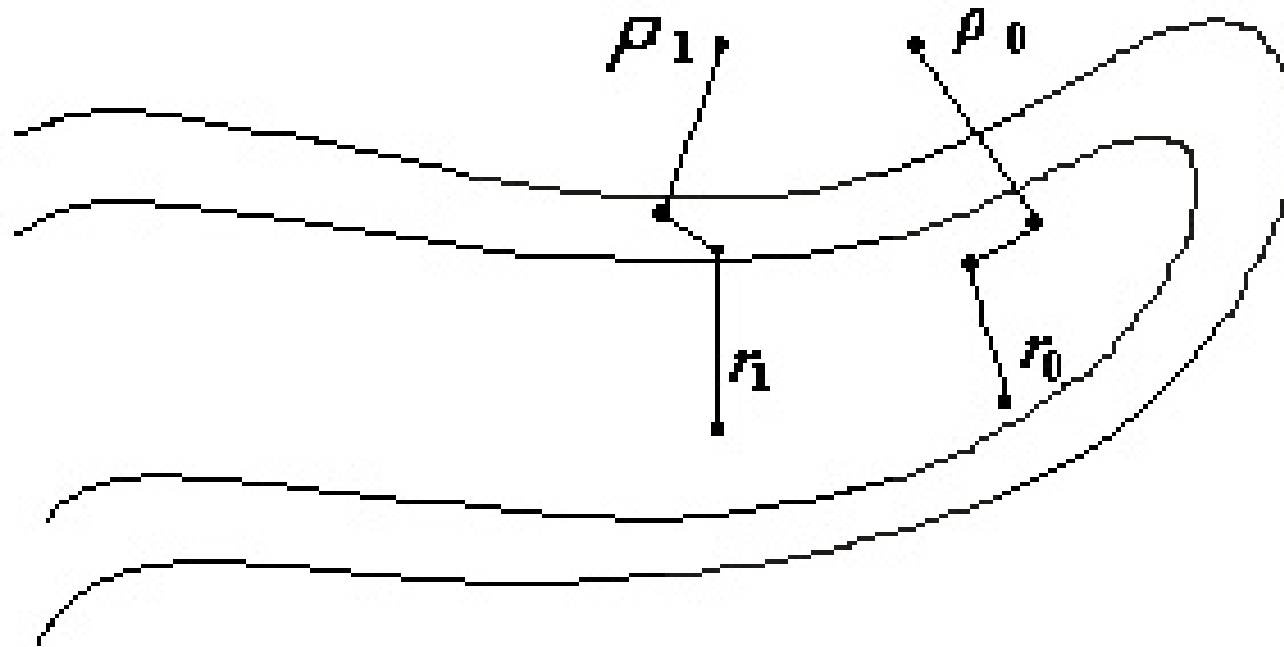


Рис. 116. Иллюстрация к методу оврагов

Проведем через точки r_0 и r_1 на дне оврага прямую – приблизительную линию дна оврага, передвинемся по этой линии в сторону убывания функции и выберем новую точку r_2 на этой прямой, на расстоянии h от точки r_1 .

Величина h называется овражным шагом и для каждой функции подбирается в ходе расчета.

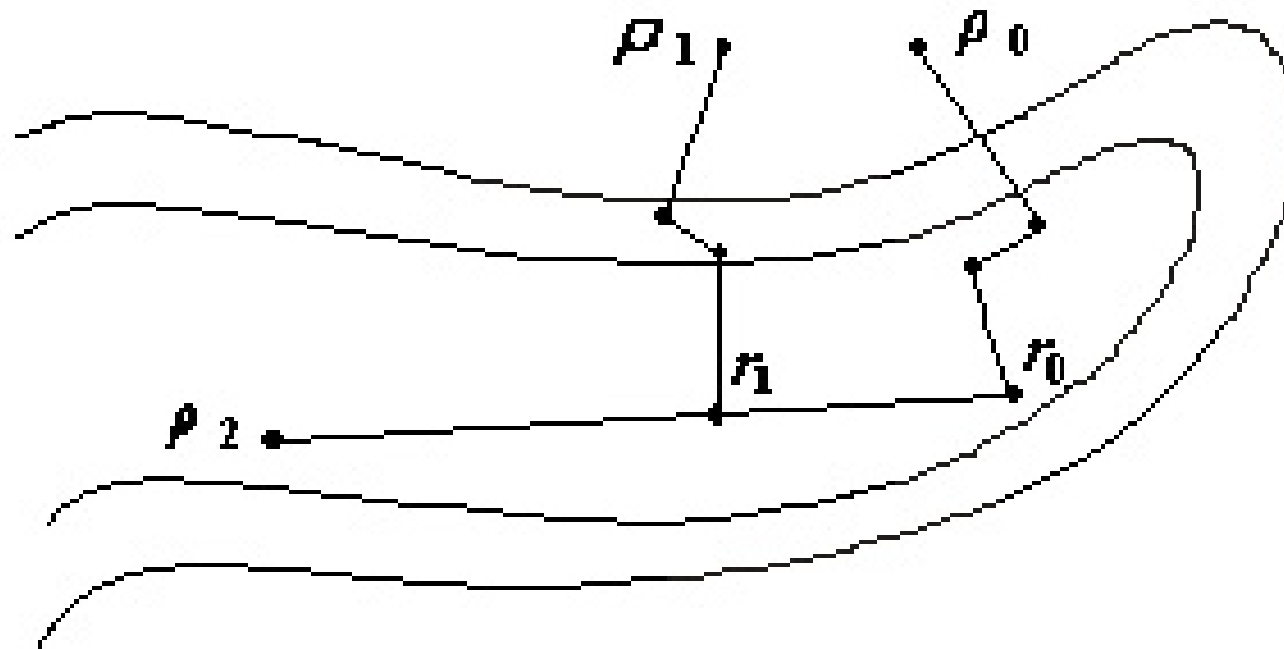


Рис. 11в. Иллюстрация к методу оврагов

Дно оврага не является отрезком прямой, поэтому точка ρ_2 на самом деле лежит не на дне оврага, а на его склоне.

Из этой точки снова спустимся на дно и попадем в некоторую точку r_2 . Затем соединим точки r_1 и r_2 прямой, наметим новую линию дна оврага и сделаем новый шаг по оврагу.

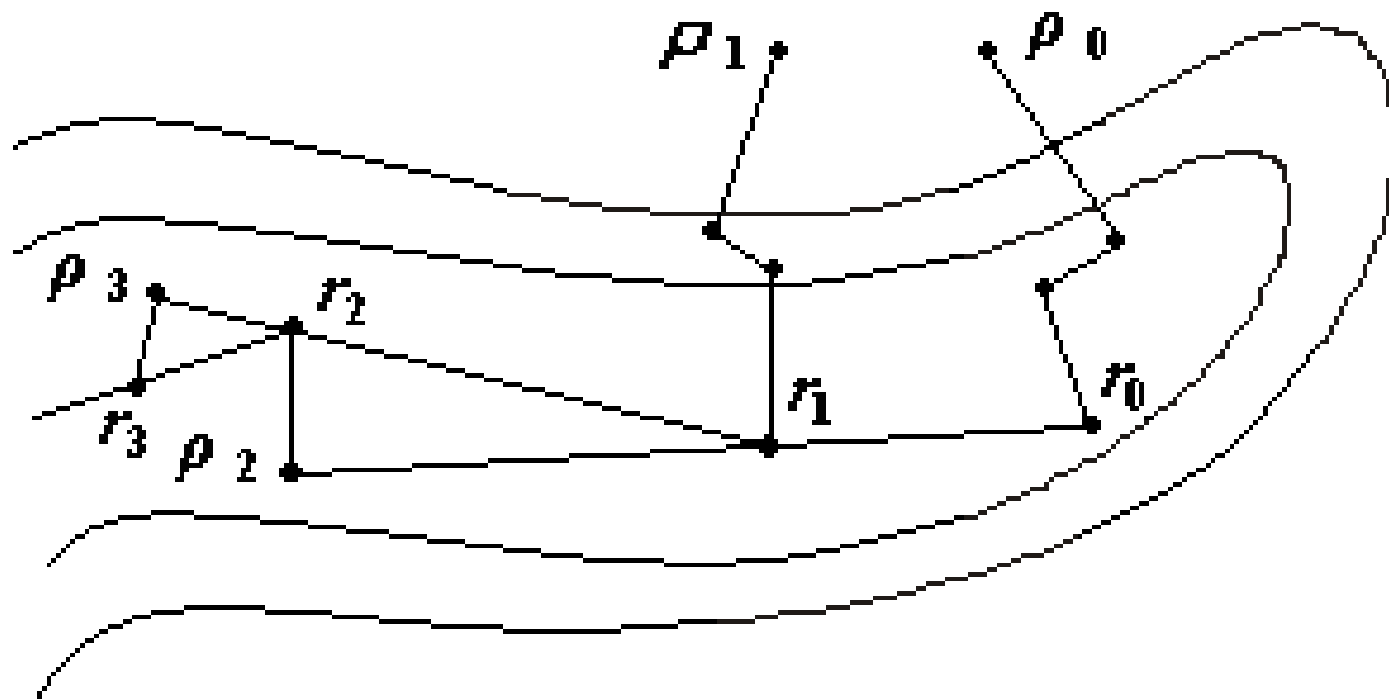


Рис. 11г. Иллюстрация к методу оврагов

Продолжим процесс до тех пор, пока значения функции на дне оврага, то есть в точках r_0, r_1, \dots, r_n убывают. В случае, когда $F(r_{n+1}) > F(r_n)$, процесс надо прекратить и точку r_{n+1} не использовать.

Метод оврагов рассчитан на то, чтобы пройти вдоль оврага и выйти в котловину около минимума.

В этой котловине значение минимума лучше уточнять другими методами.

Случайный поиск

Методы спуска не полноценны на неупорядоченном рельефе.

Если экстремумов много, то спуск из одного нулевого приближения может сойтись только к одному из локальных минимумов, не обязательно абсолютному.

Тогда для исследования задачи можно применить случайный поиск.

Предполагают, что искомый минимум лежит в некотором n -мерном параллелепипеде.

В этом параллелепипеде выбирают случайным образом N пробных точек.

Однако даже при очень большом числе пробных точек вероятность того, что хотя бы одна точка попадает в небольшую окрестность локального минимума, ничтожно мала.

Действительно, пусть $N = 10^6$ и диаметр котловины около минимума составляет 10% от пределов изменения каждой координаты.

Тогда объем этой котловины составляет 0.1^n часть объема n -мерного параллелепипеда.

Уже при числе переменных $n > 6$ практически ни одна точка в котловину не попадет.

Поэтому берут небольшое число точек $N = (5 \div 20) \cdot n$ и каждую точку рассматривают как нулевое приближение.

Из каждой точки совершают спуск, быстро попадая в ближайший овраг или котловину.

Когда шаги спуска быстро укорачиваются, его прекращают, не добиваясь высокой точности.

Этого уже достаточно, чтобы судить о величине функции в ближайшем локальном минимуме с удовлетворительной точностью. Сравнивая окончательные значения функции на всех спусках между собой, можно изучить расположение локальных минимумов и сопоставить их величины.

После этого можно отобрать нужные по смыслу задачи минимумы и провести в них дополнительные спуски для получения координат точек минимума с более высокой точностью.

МЕТОДЫ С ИСПОЛЬЗОВАНИЕМ ПРОИЗВОДНЫХ

Методы спуска и их различные модификации, стохастические методы, которые используют только значение функции, называются методами 0-го порядка.

Они обычно имеют весьма малую скорость сходимости.

Поэтому разработан ряд методов оптимизации, которые используют первые и вторые производные целевой функции (методы 1-го и 2-го порядка).

Прежде чем рассмотреть такие методы, введем ряд обозначений и напомним некоторые определения.

Вектор n -мерного пространства R^n
будем обозначать столбцом:

$$u = \begin{bmatrix} u_1 \\ \dots \\ \dots \\ u_n \end{bmatrix}$$

Будем говорить, что функция $f : R^n \rightarrow R$ непрерывно дифференцируема в точке $x \in R^n$, если производные $\frac{\partial f(x)}{\partial x_i}$, $i = 1, \dots, n$ существуют и непрерывны.

Тогда градиент функции f в точке x определяется как:

$$\nabla f(x) = \left[\frac{\partial f(x)}{\partial x_1}, \dots, \frac{\partial f(x)}{\partial x_n} \right]^T. \quad (4.1)$$

Будем говорить, что функция $f(x)$ непрерывно дифференцируема в открытой области $D \subset R^n$, если она непрерывно дифференцируема в любой точке из D .

Пусть $f : R^n \rightarrow R$ непрерывно дифференцируема на некоторой открытой выпуклой области $D \subset R^n$.

Тогда для $x \in D$ и произвольного ненулевого приращения $p \in R^n$ производная по направлению $p = [p_1, \dots, p_n]^T$ от функции $f(x)$ в точке x , определяемая как

$$\frac{\partial f(x)}{\partial p} \equiv \lim_{\varepsilon \rightarrow 0} \frac{f(x + \varepsilon p) - f(x)}{\varepsilon},$$

существует и равна $\nabla f(x)^T \cdot p$, где символом ' \cdot ' обозначено скалярное произведение.

Иначе можно записать

$$\frac{\partial f(x)}{\partial p} = \nabla f(x)^T \cdot p = \sum_{i=1}^n \frac{\partial f(x)}{\partial x_i} p_i. (4.2)$$

Будем говорить, что функция $f : R^n \rightarrow R$ дважды непрерывно дифференцируема в $x \in R^n$, если производные $\frac{\partial^2 f(x)}{\partial x_i \partial x_j}$, $1 \leq i, j \leq n$, существуют и непрерывны.

Гессианом функции (матрицей Гессе) f в точке x называется матрица размера $n \times n$, и ее (i, j) -й элемент равен

$$H_{ij} = \nabla^2 f(x)_{ij} = \frac{\partial^2 f(x)}{\partial x_i \partial x_j}, \quad 1 \leq i, j \leq n. \quad (4.3)$$

Пусть $f : R^n \rightarrow R$ дважды непрерывно дифференцируема в открытой области $D \subset R^n$.

Тогда для любого $x \in D$ и произвольного ненулевого приращения $p \in R^n$ вторая производная по направлению p от функции f в точке x , определяемая как

$$\frac{\partial^2 f(x)}{\partial p^2} \equiv \lim_{\varepsilon \rightarrow \infty} \frac{\frac{\partial f}{\partial p}(x + \varepsilon p) - \frac{\partial f}{\partial p}(x)}{\varepsilon},$$

существует и для нее выполняется равенство:

$$\frac{\partial^2 f(x)}{\partial p^2} = p^T \nabla^2 f(x) p. \quad (4.4)$$

Пусть A – действительная симметричная матрица размером $n \times n$. Будем говорить, что A *положительно определена*, если для любого ненулевого вектора $u \in R^n$ выполняется неравенство

$$u^T Au > 0.$$

Матрица A *положительно полуопределена*, если $u^T Au \geq 0$ для всех $u \in R^n$.

Для того, чтобы точка x^* была локальной точкой минимума $f(x)$ необходимо, чтобы $\nabla f(x^*) = 0$.

Достаточное условие, кроме того, требует положительной определенности $\nabla^2 f(x^*)$, а необходимое – по крайней мере положительной полуопределенности $\nabla^2 f(x^*)$

Далее будем полагать, что

$$f(x), \nabla f(x), \nabla^2 f(x)$$

существуют и непрерывны.

Описываемые ниже методы основаны на итерационной процедуре, реализуемой в соответствии с формулой

$$x^{k+1} = x^k + \lambda^k s(x^k),$$

где x^k – текущее приближение к решению x^* , λ^k – параметр, характеризующий длину шага, $s^k = s(x^k)$ – направление поиска в n -мерном пространстве.

Рассмотрим методы первого порядка, использующие первые производные.

Градиентные методы

Градиент функции в любой точке показывает направление наибольшего локального увеличения $f(x)$.

Поэтому при поиске минимума можно попробовать двигаться в направлении, противоположном градиенту в данной точке, то есть в направлении наискорейшего спуска. Такой подход приведет к итерационной формуле, описывающей

метод градиентного спуска:

$$x^{k+1} = x^k - \lambda^k \nabla f(x^k) \text{ ИЛИ}$$

$$x^{k+1} = x^k - \lambda^k \frac{\nabla f(x^k)}{\|\nabla f(x^k)\|} = x^k - \lambda^k s^k,$$

где $\|\nabla f(x^k)\|$ - норма градиента и,

соответственно, s^k - единичный вектор.

В качестве *нормы* вектора u можно выбрать так-называемую Гауссову норму

$$\| u \| = \sqrt{u_1^2 + \dots + u_n^2} .$$

В зависимости от выбора параметра λ траектория спуска будет существенно различаться.

При большом значении λ траектория будет представлять собой колебательный процесс, а при слишком больших λ процесс может расходиться.

При малых λ траектория будет плавной, но и процесс будет сходиться медленно.

Параметр λ^k можно принимать постоянным или выбирать различным на каждой итерации. Иногда на каждом k -ом шаге параметр λ^k выбирают, производя одномерную минимизацию вдоль направления s^k с помощью какого-либо одномерного метода. Обычно такой процесс называют *методом наискорейшего спуска*, или методом Коши.

Если λ^k определяется в результате одномерной минимизации, то есть $\lambda^k = \mathit{arg} \min_{\lambda} f(x^k + \lambda s^k)$, то градиент в точке очередного приближения будет ортогонален направлению предыдущего спуска (рис. 11):

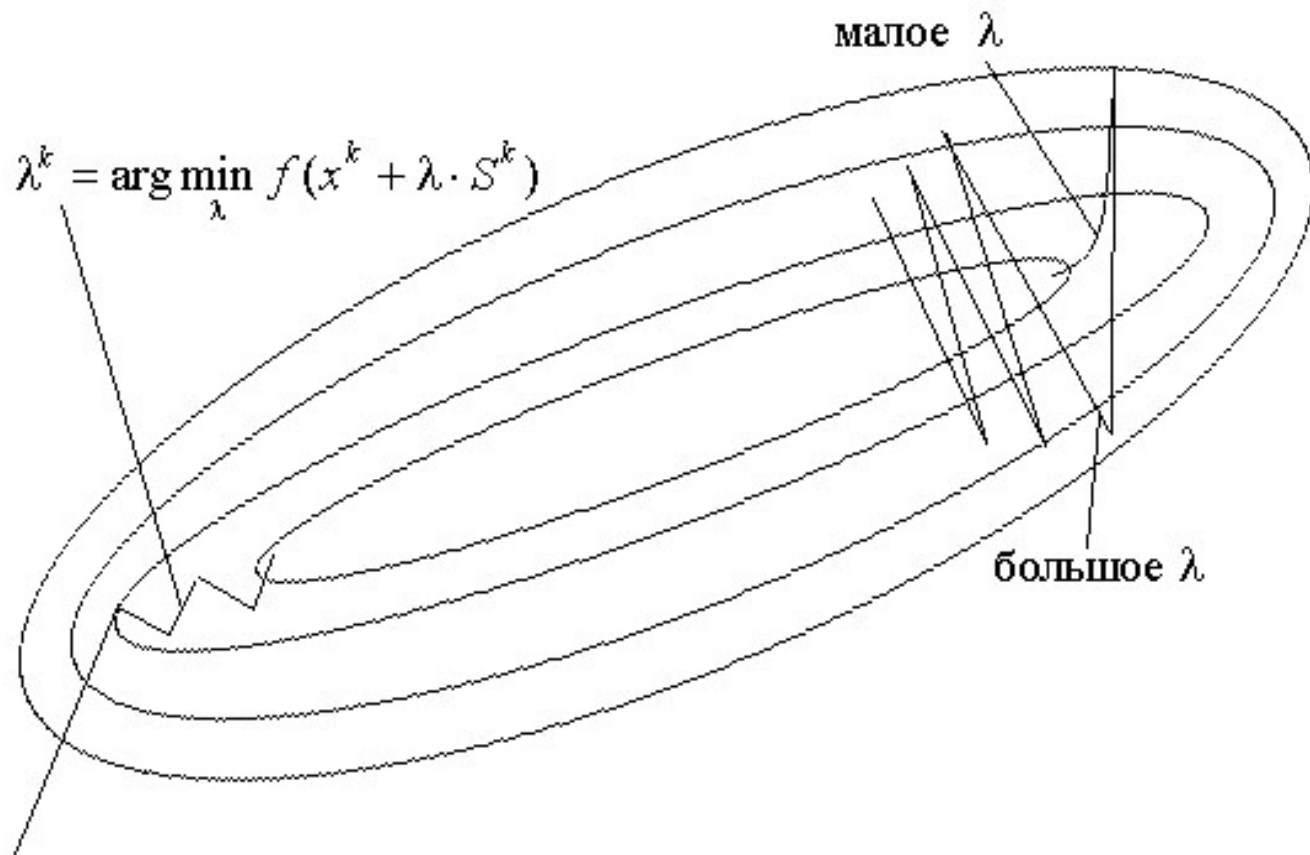


Рис. 12. Иллюстрация к градиентным методам

Одномерная оптимизация вдоль направления s^k улучшает сходимость метода, но одновременно возрастает число вычислений функции $f(x)$ на каждой итерации.

Кроме того, вблизи экстремума норма $\|\nabla f(x)\|$ близка к нулю и сходимость здесь будет очень медленной.

Эффективность алгоритма зависит от вида минимизируемой функции.

Так, для функции $f(x) = x_1^2 + x_2^2$ метод Коши сойдется к минимуму за одну итерацию при любом начальном приближении, а для функции $f(x) = x_1^2 + 100x_2^2$ сходимость будет очень медленной.

Вообще, эффективность этого метода на овражных рельефах весьма плохая. Этот метод используется очень редко.

Метод Ньютона

Построим квадратичную модель функции в окрестности точки x^k , разложив ее в ряд Тейлора до членов второго порядка:

$$\hat{f}(x^k + p) = f(x^k) + \nabla f^T(x^k)p + \frac{1}{2}p^T \nabla^2 f(x^k)p. \quad (4.5)$$

Найдем точку $x^{k+1} = x^k + s^k$ из условия минимума квадратичной модели (4.5). Необходимым условием этого минимума будет $\nabla \hat{f}(x^{k+1}) = 0$.

Имеем:

$$\nabla \hat{f}(x^k + s^k) = 0 = \nabla f(x^k) + \nabla^2 f(x^k) s^k.$$

Это приводит к следующему алгоритму:

на каждой итерации k решить систему уравнений

$$\nabla^2 f(x^k) s^k = -\nabla f(x^k) \quad (4.6)$$

относительно s^k и положить

$$x^{k+1} = x^k + s^k \quad (4.7)$$

Алгоритм (4.6), (4.7) называется методом Ньютона.

Положительной стороной этого метода является то, что если x^0 достаточно близко к точке локального минимума функции $f(x)$ с невырожденной матрицей Гессе $\nabla^2 f(x^*)$ (матрица Гессе является тогда положительно определенной), то

последовательность $\{x^k\}$, генерируемая методом Ньютона, будет сходиться к минимуму x^* и скорость сходимости будет так-называемой q -квадратичной.

К недостаткам метода относятся:

1. Метод не сходится глобально, то есть требует достаточно хорошего начального приближения x^0 ;
2. Требуется аналитического задания первых и вторых производных;

3. На каждом шаге необходимо решать систему уравнений (4.7);

4. В методе Ньютона нет ничего, что удерживало бы его от продвижения в сторону максимума или седловой точки, где градиент функции тоже равен нулю.

Здесь каждый шаг направляется просто в сторону стационарной точки локальной квадратичной модели независимо от того, является ли стационарная точка минимумом, максимумом или седловой точкой.

Этот шаг оправдан при минимизации, только когда гессиан $\nabla^2 f(x^k)$ положительно определен.

Вообще говоря, метод Ньютона не обладает высокой надежностью.

Метод Марквардта

Этот метод является комбинацией методов градиентного спуска и Ньютона, в котором удачно сочетаются положительные свойства обоих методов.

Движение в направлении антиградиента из точки, расположенной на значительном расстоянии от точки минимума, обычно приводит к существенному уменьшению целевой функции.

С другой стороны, направление эффективного поиска в окрестности точки минимума определяется по методу Ньютона.

В соответствии с методом Марквардта, направление поиска s^k определяется равенством

$$\left(H^k + \lambda^k I \right) \cdot s^k = -\nabla f \left(x^k \right), \quad (4.8),$$

а новая точка x^{k+1} задается формулой

$$x^{k+1} = x^k + s^k. \quad (4.9)$$

В системе (4.8) I – единичная матрица, H^k – матрица Гессе, $\lambda^k \geq 0$.

В формуле (4.9) коэффициент перед s^k взят равным 1, так как параметр λ^k в (4.8) позволяет менять и длину шага, и его направление.

На начальной стадии поиска параметру λ^0 приписывается некоторое большое значение, например 10^4 , так что левая часть равенства Марквардта (4.8) при больших λ^0 примет вид

$$(H^0 + \lambda^0 I)s^k \approx (\lambda^0 I)s^k = \lambda^0 s^k. \quad (4.10)$$

Таким образом, большим значениям λ^0 , как видно из (4.8) и (4.10), соответствует направление поиска

$$s^k = -\frac{1}{\lambda^0} \nabla f(x^k),$$

то есть направление наискорейшего спуска.

Из формулы

$$\left(H^k + \lambda^k I \right) \cdot s^k = -\nabla f \left(x^k \right), \quad (4.8)$$

можно заключить, что при уменьшении λ^k до нуля направление вектора s^k изменяется от противоположного градиенту, до направления, определяемому по Ньютону.

Если после первого шага получена точка с меньшим значением целевой функции, то есть $f(x^1) < f(x^0)$, следует выбрать $\lambda^1 < \lambda^0$ и реализовать еще один шаг.

В противном случае нужно положить $\lambda^0 = \beta \lambda^0$, где $\beta > 1$, и вновь реализовать предыдущий шаг.

Заметим, что недостатком метода Ньютона является то, что если матрица Гессе H^k не является положительно определенной, то Ньютоновский шаг s^k не приводит к убыванию функции.

Поэтому “исправление” Гессиана в соответствии с формулой $H^k + \lambda^k I$ модифицирует матрицу и при соответствующем выборе λ^k делает ее положительно определенной, так как единичная матрица положительно определена.

Приведем теперь алгоритм метода:

1. Задать x^0 – начальное приближение к x^* ,
 M – максимально допустимое количество итераций и ε – параметр сходимости.
2. Положить $k = 0$, $\lambda^0 = 10^4$.
3. Вычислить $\nabla f(x^k)$.
4. Проверить $\|\nabla f(x^k)\| < \varepsilon$?
Если да, то перейти к п.11.

5. Проверить $k \geq M$? Если да, то перейти к п.11.
6. Вычислить шаг s^k , решив систему
$$\left(H^k + \lambda^k I\right) \cdot s^k = -\nabla f\left(x^k\right).$$
7. Положить $x^{k+1} = x^k + s^k$.

8. Проверить: $f(x^{k+1}) < f(x^k)$?

Если да, то перейти к п.9, иначе к п.10.

9. Положить $\lambda^{k+1} = \frac{1}{2}\lambda^k$, $k = k + 1$.

Перейти к п.3.

10. Положить $\lambda^k = 2\lambda^k$. Перейти к п.6.

11. Вывод результатов:

$$x^k, f(x^k), \nabla f(x^k), \|\nabla f(x^k)\|, k.$$

Норму можно вычислять по формуле

$$\|\nabla f(x^k)\| = \sqrt{\left(\frac{\partial f(x^k)}{\partial x_1}\right)^2 + \dots + \left(\frac{\partial f(x^k)}{\partial x_n}\right)^2}$$

Отметим, что в различных модификациях метода Ньютона требуется большое количество вычислений, так как на каждой итерации следует сначала вычислить элементы матрицы $n \times n$, а затем решать систему линейных уравнений.

Применение конечно-разностной аппроксимации первых и вторых производных только ухудшит ситуацию.

УСЛОВНАЯ ОПТИМИЗАЦИЯ

Ряд инженерных задач связан с оптимизацией при наличии некоторого количества ограничений на управляемые переменные.

Такие ограничения существенно уменьшают размеры области, в которой ищется оптимум.

На первый взгляд может показаться, что уменьшение размеров допустимой области должно упростить процедуру поиска оптимума.

Однако, напротив, процесс оптимизации становится более сложным, поскольку при наличии ограничений даже нельзя использовать применяемые нами выше условия оптимальности.

При этом может нарушаться даже основное условие, в соответствии с которым оптимум должен достигаться в стационарной точке, характеризующейся нулевым градиентом.

Например, безусловный минимум функции $f(x) = (x - 2)^2$ имеет место в стационарной точке $x = 2$. Но если задача минимизации решается с учетом ограничения $x \geq 4$, то будет найден *условный минимум*, которому соответствует точка $x = 4$.

Эта точка не является стационарной точкой функции $f(x)$, так как $f'(4) = 4$.

Поэтому нужно изучить необходимые и достаточные условия оптимума в задачах с ограничениями, которые иначе называют задачами *условной оптимизации*.

Задачи с ограничениями в виде равенств

Рассмотрим задачу:

$$f(x) \rightarrow \min, x \in R^n$$

при ограничениях

$$h_k(x) = 0, k = 1, 2, \dots, K.$$

Одним из методов ее решения является метод множителей Лагранжа.

Множители Лагранжа

С помощью метода множителей Лагранжа по существу устанавливаются необходимые условия, позволяющие идентифицировать точки оптимума в задачах оптимизации с ограничениями-равенствами.

При этом задача с ограничениями преобразуется в эквивалентную задачу безусловной оптимизации, в которой фигурируют некоторые неизвестные параметры, называемые *множителями Лагранжа*.

Рассмотрим задачу с одним ограничением-равенством:

$$f(x) \rightarrow \min, x \in R^n, \quad (5.1)$$

$$h_1(x) = 0. \quad (5.2)$$

В соответствии с методом множителей Лагранжа эта задача преобразуется в следующую задачу безусловной минимизации:

$$L(x; \lambda) = f(x) - \lambda h_1(x) \rightarrow \min, x \in R^n. \quad (5.3)$$

Функция $L(x; \lambda)$ называется функцией Лагранжа.

Здесь λ – множитель Лагранжа.

Пусть при заданном значении $\lambda = \lambda^0$ безусловный минимум функции $L(x; \lambda)$ по переменной x достигается в точке $x = x^0$ и x^0 удовлетворяет уравнению

$$h_1(x^0) = 0.$$

Тогда x^0 минимизирует (5.1) с учетом (5.2), поскольку для всех значений x , удовлетворяющих (5.2),

$$h_1(x) = 0 \text{ и } \min L(x; \lambda) = \min f(x).$$

Разумеется, нужно подобрать значение $\lambda = \lambda^0$ таким образом, чтобы координата точки безусловного минимума x^0 удовлетворяла равенству (5.2). Это можно сделать, если, рассматривая λ как переменную, найти безусловный минимум функции Лагранжа (5.3) в виде функции λ , а затем выбрать значение λ , при котором выполняется равенство (5.2).

Пример.

Решить задачу

$$f(x) = x_1^2 + x_2^2 \rightarrow \min$$

при ограничении

$$h_1(x) = 2x_1 + x_2 - 2 = 0.$$

Построим функцию Лагранжа:

$$L(x; \lambda) = x_1^2 + x_2^2 - \lambda(2x_1 + x_2 - 2).$$

Определим ее безусловный минимум.

Для этого найдем стационарную точку функции Лагранжа, приравняв нулю компоненты ее градиента:

$$\frac{\partial L}{\partial x_1} = 2x_1 - 2\lambda = 0 \quad \Rightarrow \quad x_1^0 = \lambda,$$

$$\frac{\partial L}{\partial x_2} = 2x_2 - \lambda = 0 \quad \Rightarrow \quad x_2^0 = \frac{\lambda}{2}.$$

Для того чтобы проверить, соответствует ли стационарная точка x^0 минимуму, вычислим матрицу Гессе функции Лагранжа, рассматриваемой как функция от x :

$$H_L(x; \lambda) = \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix}.$$

Эта матрица положительно определена,
так как для любого ненулевого вектора
 $u^T = (a, b)$

$$\begin{aligned} u^T H_L u &= (a, b) \cdot \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix} \cdot \begin{pmatrix} a \\ b \end{pmatrix} = \\ &= (2a, 2b) \cdot \begin{pmatrix} a \\ b \end{pmatrix} = 2a^2 + 2b^2 > 0. \end{aligned}$$

Это означает, что функция Лагранжа $L(x; \lambda)$ в точке $x_1^0 = \lambda$ и $x_2^0 = \frac{\lambda}{2}$ имеет точку минимума.

Оптимальное значение λ находится путем подстановки значений x_1^0 и x_2^0 в уравнение $2x_1 + x_2 = 2$, откуда

$$2\lambda + \frac{\lambda}{2} = 2 \quad \text{и} \quad \lambda^0 = \frac{4}{5}.$$

Таким образом, условный минимум достигается при

$$x_1^0 = \frac{4}{5}, x_2^0 = \frac{2}{5},$$

а минимальное значение $f(x^0; \lambda^0)$ есть $\frac{4}{5}$.

Очень часто оказывается, что решение системы

$$\frac{\partial L}{\partial x_j} = 0, \quad j = 1, 2, \dots, n$$

в виде явной функции переменной λ получить нельзя.

Тогда значения x и λ находятся путем решения следующей системы, состоящей из $n+1$ уравнений с $n+1$ неизвестными:

$$\frac{\partial L}{\partial x_j} = 0, \quad j = 1, 2, \dots, n,$$

$$h_1(x) = 0.$$

Решить такую систему можно каким-либо численным методом.

Для каждого из решений $(x^0; \lambda^0)$ вычисляется матрица Гессе функции Лагранжа, рассматриваемой как функция от x .

Если она положительно определена, то решение – точка минимума.

Метод множителей Лагранжа можно распространить на случай, когда задача имеет несколько ограничений в виде равенств:

$$f(x) \rightarrow \min, x \in R^n,$$

$$h_k(x) = 0, k = 1, 2, \dots, K.$$

Функция Лагранжа принимает вид

$$L(x; \lambda) = f(x) - \sum_{k=1}^K \lambda_k h_k.$$

Здесь $\lambda_1, \dots, \lambda_K$ – множители Лагранжа, то есть неизвестные параметры, значения которых нужно определить. Приравнивая частные производные L по x нулю, получаем следующую систему

$$\frac{\partial L(x, \lambda)}{\partial x_1} = \frac{\partial L(x, \lambda)}{\partial x_2} = \dots = \frac{\partial L(x, \lambda)}{\partial x_n} = 0.$$

Если найти решение этой системы в виде функций от вектора λ затруднительно, то можно расширить последнюю систему путем включения в неё ограничений-равенств:

$$h_1(x) = 0, \quad h_2(x) = 0, \dots, h_K(x) = 0.$$

Решение расширенной системы, состоящей из $N+K$ уравнений с $N+K$ неизвестными, определяет стационарную точку функции L .

После этого реализуется процедура проверки на минимум или максимум, которая проводится на основе вычисления элементов матрицы Гессе функции Лагранжа, рассматриваемой как функция от x .

Методы штрафных функций

Рассмотрим задачу условной оптимизации

$$f(x) \rightarrow \min, x \in R^n, \quad (5.4)$$

$$g_j(x) \leq 0, j = 1, 2, \dots, J, \quad (5.5)$$

$$h_k(x) = 0, k = 1, 2, \dots, K, \quad (5.6)$$

$$x_i^{(l)} \leq x_i \leq x_i^{(u)}, i = 1, 2, \dots, n. \quad (5.7)$$

Такая задача также называется задачей *нелинейного программирования*.

Говорят, что точка x соответствует допустимому решению задачи нелинейного программирования, если для нее выполняются все ограничения, то есть соотношения (5.5-5.7).

Предполагается, что для вектора, являющегося решением задачи нелинейного программирования, известно некоторое начальное приближение, возможно недопустимое.

В методах штрафных функций строится последовательность точек x^m , $m = 0, 1, \dots, M$, которая начинается с заданной точки x^0 и заканчивается точкой x^M , дающей наилучшее приближение к x^* среди всех точек построенной последовательности.

В качестве x^m берутся точки решения вспомогательной задачи безусловной минимизации, полученной с помощью преобразования исходной целевой функции с помощью так называемых штрафных функций.

В этих методах исходная задача условной оптимизации преобразуется в *последовательность* задач безусловной оптимизации.

Методы штрафных функций классифицируются в соответствии со способами учета ограничений-неравенств. В зависимости от того, являются ли элементы последовательности x^m допустимыми или недопустимыми точками, говорят о *методах внутренней и внешней точки* соответственно.

Если последовательность содержит точки обоих типов, метод называют *смешанным*.

Пусть необходимо решить задачу
(5.4-5.7).

Основная идея метода штрафных функций
заключается в следующем.

Строят вспомогательную функцию

$$Q(x, r, l)$$

$$Q(x, r, l) = f(x) + \sum_{j=1}^J r_j G_j(g_j(x)) + \sum_{k=1}^K l_k H_k(h_k(x)), \quad (5.8)$$

такую, что приближенное решение задачи (5.4-5.7) получается в результате решения *последовательности* задач безусловной минимизации функции

$$Q(x, r, l) \rightarrow \min, \quad x \in R^n. \quad (5.9)$$

В методе *внешних штрафных функций* функции H , G выбираются таким образом, чтобы они становились отличными от нуля (положительными) при нарушении соответствующего ограничения.

А так как мы минимизируем (5.8), то движение в сторону нарушения становится невыгодным. Внутри допустимой области в данном методе функции H и G должны быть равны нулю.

Например, для ограничений неравенств $G_j(g_j(x)) \rightarrow 0$, при $g_j(x) \rightarrow 0^+$.

Приближенное решение задачи (5.4-5.7) получается в результате решения последовательности задач (5.9) при

$$r_j, l_k \rightarrow \infty, \quad j = 1, \dots, J, \quad k = 1, \dots, K.$$

Соответствующие методы называют *методами внешней точки*.

В методе *барьерных функций* функции H , G выбираются отличными от нуля в допустимой области и такими, чтобы при приближении к границе допустимой области (изнутри) они возрастали, препятствуя выходу при поиске за границу области.

В этом случае эти функции должны быть малыми (*положительными или отрицательными*) внутри допустимой области и большими *положительными* вблизи границы (внутри области). Например, для ограничений неравенств $G_j(g_j(x)) \rightarrow \infty$, при $g_j(x) \rightarrow 0^-$.

Такие методы называют еще *методами внутренней точки*.

В алгоритмах, использующих функции штрафа данного типа, требуют, чтобы в процессе поиска точка всегда оставалась внутренней точкой допустимой области.

Приближенное решение задачи (5.4-5.7) получается в результате решения последовательности задач (5.9) при $r_j, l_k \rightarrow 0, j = 1, \dots, J, k = 1, \dots, K.$

Для ограничений-равенств при выборе функций штрафов обычно требуют, чтобы $H_k(h_k(x)) \rightarrow 0$, при $h_k(x) \rightarrow 0$.

Это могут быть, например, функции вида:

$$H_k(h_k(x)) = |h_k(x)|, \text{ или}$$

$H_k(h_k(x)) = (h_k(x))^\alpha$, где α – четное (например, $\alpha=2$).

Для ограничений-неравенств функции штрафа подбирают таким образом, чтобы

$$G_j(g_j(x)) = 0, \text{ при } g_j(x) \leq 0,$$

$$G_j(g_j(x)) > 0, \text{ при } g_j(x) > 0.$$

Этому требованию отвечают функции вида:

$$G_j(g_j(x)) = \frac{1}{2} \{g_j(x) + |g_j(x)|\}, \quad (5.10)$$

$$G_j(g_j(x)) = \left[\frac{1}{2} \{g_j(x) + |g_j(x)|\} \right]^\alpha, \quad (5.11)$$

при четном α .

При $\alpha=2$ штраф называют квадратичным.

В качестве барьерных функций для ограничений-неравенств могут служить функции вида:

$$G_j(g_j(x)) = \frac{1}{-g_j(x)}, \quad (5.12)$$

$$G_j(g_j(x)) = -\ln(-g_j(x)). \quad (5.13)$$

Логарифмический штраф (5.13) – это барьерная функция, не определенная в недопустимых точках (то есть для таких x , что $g(x) > 0$).

Поэтому в тех случаях, когда приходится иметь дело с недопустимыми точками (например, когда заданное начальное приближение не является допустимым), требуется специальная процедура, обеспечивающая попадание в допустимую область.

Штраф, заданный функцией (5.12), не имеет отрицательных значений в допустимой области.

Этот штраф, как и предыдущий, является барьером; в этом случае также возникают трудности, связанные с возможным появлением недопустимых точек.

Часто функцию $Q(x, r, l)$ выбирают в виде

$$Q(x, r, l) = f(x) - R \sum_{j=1}^J \frac{1}{g_j(x)} + \frac{1}{R} \sum_{k=1}^K (h_k(x))^2,$$

ИЛИ

$$Q(x, r, l) = f(x) - R \sum_{j=1}^J \ln(-g_j(x)) + \frac{1}{R} \sum_{k=1}^K (h_k(x))^2,$$

где положительный штрафной параметр $R \rightarrow 0$, монотонно убывая от итерации к итерации.

Последовательность действий при реализации методов штрафных функций или барьерных функций выглядит следующим образом:

1. На основании задачи (5.4-5.7) строим функцию (5.8). Выбираем начальное приближение x и начальные значения коэффициентов штрафа r_j, l_k , число итераций, точность безусловной оптимизации, точность соблюдения ограничений и т.д.

2. Решаем задачу (5.9).

3. Если полученное решение не удовлетворяет системе ограничений в случае использования метода штрафных функций, то увеличиваем значения коэффициентов штрафа и снова решаем задачу (5.9).

В случае метода барьерных функций значения коэффициентов уменьшаются, чтобы можно было получить решение на границе.

4. Процесс прекращается, если найденное решение удовлетворяет системе ограничений с определенной точностью.

Метод факторов

Своеобразным и очень эффективным методом штрафов является метод факторов (или множителей), который основан на штрафе типа “квадрат срезки” для ограничений-неравенств.

Такой штраф определяется следующим образом:

$$S = R \cdot \langle g(x) \rangle^2, \quad (5.14)$$

где *срезка* t определяется так:

$$\langle t \rangle = \begin{cases} t, & \text{если } t \geq 0, \\ 0, & \text{если } t < 0. \end{cases} \quad (5.15)$$

Этот штраф внешний и стационарные точки функции $Q(x, R)$ могут оказаться недопустимыми.

С другой стороны, недопустимые точки не создают в данном случае дополнительных сложностей по сравнению с допустимыми. Различие между ними состоит лишь в том, что в допустимых точках штраф равен нулю.

В методе факторов на каждой итерации производится безусловная минимизация функции

$$\begin{aligned} Q(x, \sigma, \tau) = & f(x) + \\ & + R \sum_{j=1}^J \left\{ \left\langle g_j(x) + \sigma_j \right\rangle^2 - \sigma_j^2 \right\} + \\ & + R \sum_{k=1}^K \left\{ \left[h_k(x) + \tau_k \right]^2 - \tau_k^2 \right\}, \end{aligned} \quad (5.16)$$

где R – постоянный весовой коэффициент, а угловые скобки обозначают *операцию срезки*.

Параметры (факторы) σ и τ осуществляют сдвиг штрафных слагаемых.

Компоненты векторов σ и τ меняются по ходу вычислений, однако в процессе решения каждой вспомогательной безусловной задачи оба эти вектора остаются постоянными.

Начальные значения факторов σ и τ можно выбрать нулевыми.

Обозначим через x^m точку минимума функции $Q(x, \sigma^m, \tau^m)$, используемой на m -ой итерации.

При переходе к $(m+1)$ -й итерации факторы пересчитываются по формулам

$$\sigma_j^{m+1} = \left\langle g_j \left(x^m \right) + \sigma_j^m \right\rangle, \quad j = 1, \dots, J, \quad (5.17)$$

$$\tau_k^{m+1} = h_k \left(x^m \right) + \tau_k^m, \quad k = 1, \dots, K. \quad (5.18)$$

Формулы пересчета таковы, что в результате сдвига при переходе к новой подзадаче штраф за нарушение ограничений возрастает, и вследствие этого точки x^m приближаются к допустимой области.

Для контроля сходимости метода
используют последовательности

$$x^m, \sigma^m, \tau^m, f(x^m), g(x^m), h(x^m).$$

Прекращение основного процесса происходит, когда члены, по крайней мере, одной из этих последовательностей, перестают изменяться при пересчете факторов и последующей безусловной минимизации.

Заметим, что величина положительного параметра R влияет на свойства метода, но конструктивного алгоритма его выбора не существует.