

- Вблизи точки экстремума  $M^*$  сходимость координатного спуска и по координатам, и по градиенту *линейная* (достаточно медленная, что с практической точки зрения плохо);
- по "циклам" спусков можно делать ускорения по методу Эйткена;
- При попадании траектории спуска в разрешимый овраг расчет практически невозможен (слишком медленная сходимость при произвольной ориентации оврага относительно координатных осей). Поэтому выгоднее использовать методы, обладающие повышенным порядком точности.

### 3.5 Градиентные методы минимизации

В общем случае для траектории спуска  $\{M_k\}$  :  $\Phi_{k+1} < \Phi_k$  при минимизации достаточно гладких функций можно сформулировать *достаточные* условия сходимости соответствующего метода спуска, характеризующие изменение функции  $\Phi$  и её градиента  $\vec{g} = \text{grad}\Phi$  на траектории  $\{M_k\}$ .

Пусть очередной шаг совершается вдоль направления  $\vec{p}_k$  и приводит нас в точку  $M_{k+1}$ :

$$\vec{x}_{k+1} = \vec{x}_k + \vec{p}_k h_k.$$

Шаг  $h_k$  выбирается из условия минимальности  $\Phi(M)$  вдоль  $\vec{p}_k$

$$h_k : \varphi(h_k) = \min_h \varphi(h) = \min_h \Phi(\vec{x}_k + h\vec{p}_k).$$

Сформулируем достаточные условия сходимости метода спуска.

**Теорема 2.** Пусть

- 1)  $\Phi(\vec{x})$  – дважды дифференцируемая функция;
- 2) множество уровня

$$D(\Phi(\vec{x}_0)) = \{\vec{x} : \Phi(\vec{x}) \leq \Phi(\vec{x}_0)\}$$

ограничено и замкнуто;

- 3) на каждой итерации

- а) направление  $\vec{p}_k$  – "существенное направление спуска":

$$\exists \beta < 0, \quad \vec{p}_k \vec{g}_k \leq \beta < 0$$

- б)  $\Phi(x)$  "существенно убывает" (т.е. выбрано соответствующее ограничение на шаг):

$$\begin{aligned} \exists \mu_1, \quad \mu_2 : 0 < \mu_1 \leq \mu_2 \leq 1 \\ -\mu_1 h_k \vec{g}_k \cdot \vec{p}_k \leq \Phi_k - \Phi_{k+1} \leq -\mu_2 h_k \underbrace{\vec{g}_k \cdot \vec{p}_k}_{\text{отнц. число}} \end{aligned}$$

Тогда

$$\lim_{k \rightarrow \infty} \|\vec{g}_k\| = 0; \quad (M_k \rightarrow M^*)$$

т.е. метод спуска обладает сходимостью (как правило — линейной).

В основном соответствующие методы спуска отличаются выбором очередного направления  $\vec{p}_k$  и шага  $h_k$ :

**Метод "наискорейшего" спуска.** Рассмотрим линейную аппроксимацию целевой функции  $\Phi(\vec{x})$  в окрестности точки  $\vec{x}_k$ . Опираясь на формулу Тейлора:

$$\Phi(\vec{x}_k + \vec{p}) = \Phi(\vec{x}_k) + (\text{grad}\Phi(\vec{x}_k), \vec{p}) + o(\|\vec{p}\|),$$

с определенной точки зрения (локально!) естественно искать направление, по которому  $\frac{\partial \Phi}{\partial p} \equiv \vec{g}_k \cdot \vec{p}$  наибольшее по модулю отрицательное число. Это направление в первом порядке по  $\|\vec{p}\|$  обеспечивает наибольшее убывание функции  $\Phi$ .

Итак, необходимо найти направление  $\vec{p}$

$$\begin{cases} \min(\vec{g}_k \cdot \vec{p}) \\ \|\vec{p}\| = 1 \end{cases} \quad \text{— задача на} \\ \text{условный} \\ \text{экстремум} \\ \text{для } \vec{p}$$

Решение полученной задачи зависит от вида рассматриваемой нормы. Если выбрать  $C$ -энергетическую норму  $\|\vec{p}\|^2 = (C\vec{p}, \vec{p})$ , где  $C > 0$  и симметрична, тогда направление  $\vec{p}$  (с точностью до нормировочной  $Const$ )

$$\vec{p} = -C^{-1} \cdot \vec{g}_k. \quad *1)$$

Для евклидовой нормы —  $C \equiv E$  и  $p = -\vec{g}_k$ , что приводит нас к *методу наискорейшего спуска*.

$$\begin{cases} \vec{x}_{k+1} = \vec{x}_k - h_k \vec{g}_k \\ h_k : \varphi(h_k) = \min_h \Phi(\vec{x}_k - h \vec{g}_k) \end{cases} \quad (16)$$

**Замечания:**

- 1) При таком выборе  $\vec{p}_k$  и  $h_k$  (16) траектория спуска перпендикулярна линии уровня  $\Phi(x_k)$  в точке  $x_k$ .
- 2) По сходимости *наискорейший спуск* не лучше, чем координатный спуск, т.е. он обладает лишь линейной сходимостью.
- 3) Анализ сходимости наискорейшего спуска на квадратичной функции с симметричной и положительно определенной матрицей (что характерно для гессиана в окрестности невырожденного минимума)

$$\Psi(x) = \frac{1}{2}(Ax, x) + (\vec{b}, x) + C : A > 0, A^T = A$$

\*1) Показать!

дает лишь линейную сходимость. Поскольку  $A > 0$ ,  $A^T = A$  следовательно все собственные значения матрицы  $A$  положительны  $\forall i \lambda_i(A) > 0$ . Сходимость метода наискорейшего спуска характеризуют величиной

$$\varkappa = \frac{\lambda_{max}(A)}{\lambda_{min}(A)} = \|A\| \cdot \|A^{-1}\| = CondA$$

$$\Psi(x_{k+1}) - \Psi(x^*) \simeq \left( \frac{\varkappa - 1}{\varkappa + 1} \right)^2 (\Psi(\vec{x}_k) - \Psi(x^*)). \quad (17)$$

Полученная оценка скорости сходимости, например для  $\varkappa = 100$  (хорошая обусловленность матрицы  $A$ ) даёт  $q \approx 0,96(!)$  и нужны сотни итераций для уменьшения погрешности на порядок.

Расчетные формулы наискорейшего спуска (16) в этом случае принимают вид:

$$\begin{aligned} \vec{g} &= A\vec{x} + \vec{b}; \text{ Hess}\Psi = A \Rightarrow \vec{p}_k = -\vec{g}_k, \\ \psi(h) &= \Psi(\vec{x} + h\vec{p}_k) = \Psi(\vec{x}) + h(A\vec{x} + \vec{b}, \vec{p}_k) + \frac{h^2}{2}(A\vec{p}_k, \vec{p}_k), \\ \frac{\partial \psi}{\partial h} = 0 &\Leftrightarrow h_k = \left\{ \begin{array}{c} \text{получить} \\ \text{самостоятельно} \\ \text{расчетные} \\ \text{формулы} \end{array} \right\}. \end{aligned} \quad (18)$$

Тем не менее:

- 1) Необходимо бесконечное число итераций для нахождения экстремума даже в случае квадратичной функции.
- 2) Метод наискорейшего спуска не рекомендуется как серьезная минимизационная процедура. Дело в том, что свойство наискорейшего спуска является лишь *локальным* свойством, поэтому необходима частая смена направлений спуска и относительно малый шаг движения по каждому направлению, что и приводит в итоге к неэффективной вычислительной процедуре (например в случае разрешимого оврага).
- 3) Метод наискорейшего спуска невозможно адаптировать для использования информации о вторых производных  $\Phi(\vec{x})$ .

### 3.6 Методы второго порядка

**Ньютоновские методы.** Эта группа методов основана на более точной аппроксимации целевой функции в окрестности точки  $\vec{x}_k$

$$\Phi(\vec{x}_k + \vec{p}) = \underbrace{\Phi(\vec{x}_k) + \vec{g}_k \cdot \vec{p} + \frac{1}{2}(G_k \vec{p}, \vec{p})}_{\Psi(\vec{p})} + o(\|\vec{p}\|^2).$$

Минимизируемая функция  $\Psi(\vec{p})$ . Соответствующее направление и шаг берут из условия минимума  $\Psi(\vec{p})$ :

$$\left. \begin{aligned} \text{grad}\Psi = 0 \quad \Leftrightarrow \quad G_k \vec{p} + \vec{g}_k = 0; \quad \Leftrightarrow \quad \underbrace{\vec{p}_k = -G_k^{-1} \cdot \vec{g}_k}_{\text{Ньютоновское направление}} \\ \vec{x}_{k+1} = \vec{x}_k + \vec{p}_k = x_k - G_k^{-1} \cdot \vec{g}_k \end{aligned} \right\} \quad (19)$$

- Для квадратичной целевой функции  $\Psi(\vec{p})$  метод (19) решает задачу минимизации за одну(!) итерацию.
- В окрестности невырожденного экстремума имеет *квадратичную* сходимость (гессиан  $G_k > 0$  и симметричен).
- Ньютоновское направление – это направление *наискорейшего* спуска в  $G$ -энергетической метрике

$$\|\vec{p}\| = \sqrt{(G\vec{p}, \vec{p})}.$$

- Существенным является то, что на каждом шаге необходимо решать систему линейных уравнений (19) для определения *ньютоновского направления* очередной итерации.
- При модификации метода Ньютона, когда гессиан фиксируется на определенное число итераций  $G_{k_0}$  — в методе Ньютона-Рафсона — существует алгоритмический выигрыш, но при этом обеспечена лишь линейная сходимость метода.

**Метод сопряженных градиентов.** Методы *координатного спуска* или *наискорейшего спуска* требовали даже для минимизации квадратичной функции бесконечного числа итераций.

Опираясь на тейлоровское разложение в окрестности невырожденного экстремума  $x^*$  выгодно строить методы спуска, которые, по крайней мере, эффективны для квадратичных функций.

Таковыми методами, не требующими решения СЛАУ (19) на каждом итерационном шаге для определения направления спуска, являются методы *сопряженных направлений*.

Для квадратичной функции  $\Psi(\vec{x})$ :

$$\Psi(x) = \frac{1}{2}(Ax, x) + (b, x) + c, \quad A > 0, \quad A^T = A$$

они позволяют не более чем за  $n$  шагов спуска получить её минимум. Напомним:

Симметричная положительноопределенная матрица  $A > 0$ ,  $A^T = A$  – позволяет ввести "A-энергетическую" норму вектора

$$\|x\|_A = \sqrt{(Ax, x)}$$

и соответствующее скалярное произведение

$$(x, y)_A = (Ax, y) = (x, Ay).$$

**Определение** Векторы, ортогональные в  $A$ -энергетическом смысле, называются сопряженными относительно матрицы  $A$ .

$$x \underset{A}{\perp} y \Leftrightarrow (x, y)_A = (Ax, y) = (x, Ay) = 0.$$

Сопряженные векторы обладают рядом "хороших" свойств:

- 1) Если  $\{x_i\}_k$  – система сопряженных векторов и  $k \leq n$ , то эта система векторов – линейно независима.

Действительно, пусть  $\vec{x}_1 = \sum_{i=2}^k \alpha_i \vec{x}_i$  – ненулевая комбинация остальных векторов. Тогда

$$(x_1, Ax_1) = (x_1, A \sum_{i=2}^k \alpha_i \vec{x}_i) = \sum_{i=2}^k \alpha_i (x_1, Ax_i) \equiv 0$$

но  $A > 0$  и следовательно  $\vec{x}_1$  нулевой вектор, что невозможно ■

- 2) Если число векторов в рассматриваемой системе  $k = n$ , то  $\{x_i\}_n$  – сопряженный базис. Можно считать его сопряженным ОНБ, т.е.  $(x_i, x_j)_A = \delta_{ij}$ . Разложим направление  $\vec{p}$  по ОНБ  $\{x_i\}_n$  и рассмотрим квадратичную функцию на этом направлении

$$\begin{aligned} \Psi(\vec{x} + \vec{p}) &= \Psi(\vec{x}) + (Ax + b, \vec{p}) + \frac{1}{2}(A\vec{p}, \vec{p}) = \left| p = \sum \alpha_i \vec{x}_i \right| = \\ &= \Psi(\vec{x}) + (Ax + b, \sum_i \alpha_i \vec{x}_i) + \frac{1}{2} \left( A \sum_i \alpha_i \vec{x}_i, \sum_k \alpha_k \vec{x}_k \right) = \\ &= \underbrace{\sum_i \left\{ \frac{1}{2} \alpha_i^2 + \alpha_i (Ax + b, x_i) \right\}}_{n \text{ независимых слагаемых}} + \Psi(\vec{x}); \end{aligned} \quad (*)$$

Движение по каждому из сопряженных направлений  $x_i$  изменяет только одно слагаемое в сумме (\*) и, тем самым, за не более, чем  $n$  шагов приводит к минимуму функции  $\Psi$ .

Существуют различные способы построения сопряженных относительно  $A$  направлений, в частности – метод *сопряженных градиентов* (метод Флетчера-Ривса) – приводит к одной из наиболее эффективных процедур многомерной численной минимизации.

Рассмотрим снова квадратичную аппроксимацию  $\Psi(x)$  целевой функции  $\Phi(x)$  в окрестности точки  $\vec{x}_k$ :

$$\Phi(\vec{x}_k + \vec{p}) = \underbrace{\Phi(x_k) + (grad\Phi(x_k), \vec{p}) + \frac{1}{2}(hess\Phi(x_k)\vec{p}, \vec{p})}_{\Psi_k(\vec{p})} + o(\|\vec{p}\|^2).$$

На каждом *цикле* итерационных шагов для построения *сопряженного базиса* будем использовать одну и ту же матрицу  $G_k \equiv hess\Phi(x_k)$ . При этом мы будем считать, что находимся в достаточно малой окрестности точки минимума  $x^*$ , где  $G(x_k) > 0$ .

В методе сопряженных градиентов совокупность сопряженных относительно  $G \equiv G(x_k)$  направлений строится следующим образом. Опишем процедуру построения одного цикла минимизации, содержащего  $n$  шагов и точно минимизирующего  $\Psi_k(\vec{p})$ .

$$\begin{array}{l}
 \text{Цикл} \\
 \text{движения} \\
 \\
 \text{1-ый} \\
 \text{шаг:} \\
 \\
 \text{2-ой} \\
 \text{шаг:}
 \end{array}
 \quad
 \begin{array}{l}
 M_k \equiv M_k \xrightarrow{\vec{p}_1^{(1)}} M_k \xrightarrow{\vec{p}_2^{(2)}} \dots \xrightarrow{\vec{p}_{n-1}^{(n-1)}} M_k \xrightarrow{\vec{p}_n^{(n)}} M_{k+1} \\
 \\
 \vec{p}_1 = -\vec{g}_1; \quad x_k = x_k + h_1 \vec{p}_1; \quad h_1 : \psi(h_1) = \min_h \Psi_k \left( x_k + h \vec{p}_1 \right); \\
 \\
 \vec{p}_2 = -\vec{g}_2 + \alpha_1 \vec{p}_1; \quad \alpha_1 = \frac{(g_2, g_2)}{(g_1, g_1)}; \quad \vec{p}_2 \perp \vec{p}_1 \text{ отн-но } G_k
 \end{array}
 \tag{20}$$

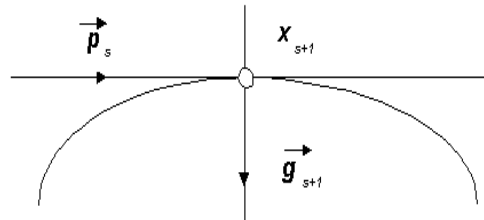
Пусть  $\vec{p}_1, \dots, \vec{p}_s$   $G_k$ -сопряженная система векторов

$$\left. \begin{array}{l}
 \vec{p}_{s+1} = -\vec{g}_{s+1} + \alpha_s \vec{p}_s \quad (\text{все остальные } \alpha_i = 0) \\
 \alpha_s = \frac{g_{s+1}^2}{g_s^2}; \quad (\text{из сообр. } (\vec{p}_{s+1}, \vec{p}_s)_{G_k} = 0) \\
 \vec{x}_{s+2} = \vec{x}_{s+1} + h_{s+1} \cdot \vec{p}_{s+1} \\
 h_{s+1} : \psi(h_{s+1}) = \min_h \Psi_k(\vec{x}_{s+1} + h \vec{p}_{s+1})
 \end{array} \right\}$$

Покажем, что (20) определяет систему сопряженных относительно  $G_k$  векторов движения  $\{\vec{p}_s\}_n$ .

а) Проверить самостоятельно 2-й шаг;

б) 1:  $\vec{g}_{s+1}$  ортогонально всем предыдущим  $\vec{p}_j$  при  $j \leq s$ , ибо спускаясь на предыдущем,  $S$ -ом шаге, мы пришли в точку  $\vec{x}_{s+1} = \vec{x}_s + h_s \vec{p}_s$  вдоль направления  $\vec{p}_s$ .



Но эта точка —  $\vec{x}_{s+1}$  — точка "минимума", т.е.  $\vec{g}_{s+1} \perp \vec{p}_s$ ,  $(\vec{g}_{s+1}, \vec{p}_s) = 0$ . Если проследить "вглубь" траектории, то

$$\vec{x}_{s+1} = \vec{x}_s + h_s \vec{p}_s = \vec{x}_{s-1} + h_{s-1} \vec{p}_{s-1} + h_s \vec{p}_s = \dots = \vec{x}_{j+1} + \sum_{j+1}^s h_i \vec{p}_i, \quad 1 \leq j \leq S-1.$$

Тогда

$$Gx_{s+1} = Gx_{j+1} + \sum_{j+1}^s h_i G \vec{p}_i.$$

Добавим слева и справа по  $\vec{b} \equiv \vec{g}(M_k)$ , и учтём, что  $G_k \equiv G$ ;  $Gx + b \equiv \vec{g}(x)$ . Таким образом

$$\vec{g}_{S+1} = \vec{g}_{j+1} + \sum_{j+1}^S h_i G p_i.$$

Тогда

$$\begin{aligned} \vec{g}_{S+1} \cdot \vec{p}_j &= \underbrace{\vec{g}_{j+1} \cdot \vec{p}_j}_{=0 \text{ для этого шага}} + \sum_{j+1}^S h_i \cdot \underbrace{(G p_i, p_j)}_{=0 \text{ в силу индукции}} \Rightarrow \\ &\Rightarrow (\vec{g}_{S+1}, \vec{p}_j) = 0, \quad 1 \leq j \leq S-1; S. \end{aligned}$$

2: Покажем, что вектор  $\vec{g}_{S+1}$  ортогонален всем градиентам  $\vec{g}_j, j = \overline{1, S}$ . Имеем

$$\vec{p}_j = -g_j + \alpha_{j-1} \vec{p}_{j-1} \Leftrightarrow \underbrace{\vec{g}_{S+1} \cdot \vec{p}_j}_{=0} = -(g_{S+1}, g_j) + \underbrace{\alpha_{j-1} \cdot (\vec{g}_{S+1}, \vec{p}_{j-1})}_{=0}$$

т.о.

$$(\vec{g}_{S+1}, \vec{g}_j) = 0, \quad j = \overline{1, S}.$$

3: Рассмотрим очередное направление:

$$\vec{p}_{S+1} = -\vec{g}_{S+1} + \alpha_S \vec{p}_S; \quad \alpha_S = \frac{g_{S+1}^2}{g_S^2}$$

и покажем, что  $\vec{p}_{S+1}$  сопряжено всем  $\vec{p}_j, j \leq S$ . Оно сопряжено, по крайней мере, со всеми  $\vec{p}_j$  до предыдущего, т.е.  $(\vec{p}_{S+1}, \vec{p}_j)_{G_k} = 0, j = \overline{1, S-1}$ . Действительно, поскольку  $j \leq S-1$ , то

$$\begin{aligned} (\vec{p}_{S+1}, G \vec{p}_j^*) &= \left( -\vec{g}_{S+1} + \alpha_S \overbrace{\vec{p}_S^*, G \vec{p}_j^*}^{\text{сопряжены}} \right) = - \left( g_{S+1}, G \frac{x_{j+1} - x_j}{h_j} \right) = \\ &= - \left( g_{S+1}, \frac{(G x_{j+1} + b_k) - (G x_j + b_k)}{h_j} \right) = - \left( g_{S+1}, \frac{g_{j+1} - g_j}{h_j} \right) \equiv 0. \end{aligned}$$

Предыдущее направление:

$$\begin{aligned} (\vec{p}_{S+1}, G \vec{p}_S^*) &= -(g_{S+1}, G p_S) + \alpha_S (p_S, G p_S) = \\ &= - \left( g_{S+1}, G \frac{x_{S+1} - x_S}{h_S} \right) + \alpha_S \left( -g_S + \alpha_{S-1} \vec{p}_{S-1}, \frac{g_{S+1} - g_S}{h_S} \right) = \\ &= - \frac{g_{S+1}^2}{h_S} + \frac{g_{S+1}^2}{h_S^2} \frac{h_S^2}{h_S} = 0 \blacksquare \end{aligned}$$

Метод Флетчера-Ривса обладает квадратичной сходимостью в достаточно малой окрестности точки  $\vec{x}^*$ . Рестарт в точке  $M_k$  осуществляется по антиградиенту  $(-\vec{g}_k)$ .

Это один из наиболее эффективных методов численной минимизации функций многих переменных.

## §4. Задача минимизации функционала

### 4.1 Постановка задачи

Если любому  $y(x) \in Y$  поставлено в соответствие число  $\Phi[y(x)]$ , то говорят, что на множестве функций  $Y$  задан функционал  $\Phi[y(x)]$ .

В задаче минимизации функционала требуется: *найти  $y^*(x) \in Y$ , на которой функционал достигает своей точной нижней грани (абсолютный экстремум)*

$$y^* : \quad \Phi^* \equiv \Phi[y^*(x)] = \inf_Y \Phi[y(x)]. \quad (21)$$

В такой постановке (21) называется задачей *минимизации по аргументу*, в отличие от

$$\Phi^* \equiv \inf_Y \Phi[y(x)] \quad (21')$$

задачи *минимизации значений* функционала.

Не всякий функционал и не на всяком множестве имеет минимум. Скажем, если функционал неограничен снизу на заданном множестве, или соответствующее множество некомпактно в себе, или если функционал разрывен и т.д. Мы не будем исследовать постановку задачи (21), а будем предполагать, что (21) поставлена корректно, то есть её решение  $y^*(x)$  на  $Y$  существует, единственно и устойчиво относительно малых возмущений входных данных.

Постановка задачи (21) возникает, как правило, когда сама модель сформулирована соответствующим образом, например рассматривается функционал "действия" (или нечто похожее)

$$\Phi[y(x)] = \int_a^b F(x, y, y', \dots, y^{(p)}) dx. \quad (*)$$

Обычно к задаче (21) приводит использование вариационных методов решения "операторного" уравнения  $A[y(x)] = f(x)$ . Рассмотренный нами *метод наименьших квадратов* даёт задачу минимизации функционала "невязки" для этого уравнения

$$\Phi[y(x)] = \| Ay - f \|^2 = \int_a^b (Ay(x) - f(x))^2 \rho(x) dx, \quad \rho(x) > 0.$$

В случае *некорректно* поставленной задачи  $A[y(x)] = f(x)$  её *регуляризация* приводит к задаче минимизации *сглаживающего функционала Тихонова*

$$M_\alpha[y(x)] = \| Ay - f \|^2 + \alpha \Omega[y(x)],$$

где  $\Omega[y(x)]$  — функционал со свойствами нормы (то есть с его помощью на  $Y$  вводится структура нормированного пространства и множество  $\Omega[y] \leq Const$  компактно в  $Y$  в введенной метрике). Тогда решение задачи минимизации  $y_\alpha(x)$  при определённом способе согласования параметра регуляризации  $\alpha$  с априорной информацией о