Метод Ньютона. Квазиньютоновские методы. Матрица предобработки Методы оптимизации

Александр Безносиков

Московский физико-технический институт

5 октября 2023



• Рассмотрим задачу поиска «корня» функции:

Найти
$$t^*$$
, что $\varphi(t^*)=0$,

где
$$\varphi:\mathbb{R} o \mathbb{R}$$
.

Рассмотрим задачу поиска «корня» функции:

Найти
$$t^*$$
, что $\varphi(t^*)=0$,

где $\varphi:\mathbb{R} \to \mathbb{R}$.

ullet Пусть мы находимся в точке t^0 и хотим найти такую поправку Δt , что $t^0 + \Delta t_{\rm i} \approx t^*$.

$$\varphi(t^{\circ} + \Delta t) = \varphi(t^{\circ}) + \varphi'(t^{\circ}) \Delta t + o(\Delta t)$$

$$\approx \varphi(t^{\circ}) + \varphi'(t^{\circ}) \Delta t = 0$$

$$\varphi(t^{\circ}) + \varphi'(t^{\circ}) \Delta t = 0$$

$$\psi(t^{\circ}) + \varphi'(t^{\circ}) \Delta t = 0$$

• Рассмотрим задачу поиска «корня» функции:

Найти
$$t^*$$
, что $\varphi(t^*)=0$,

где $\varphi:\mathbb{R} \to \mathbb{R}$.

- ullet Пусть мы находимся в точке t^0 и хотим найти такую поправку Δt , что $t^0 + \Delta t pprox t^*$.
- Разложим в ряд:

$$\varphi(t^0 + \Delta t) = \varphi(t^0) + \varphi'(t^0)\Delta t + o(\Delta t).$$

• Рассмотрим задачу поиска «корня» функции:

Найти
$$t^*$$
, что $\varphi(t^*)=0$,

где $\varphi:\mathbb{R}\to\mathbb{R}$.

- ullet Пусть мы находимся в точке t^0 и хотим найти такую поправку Δt , что $t^0 + \Delta t pprox t^*$.
- Разложим в ряд:

$$\varphi(t^0 + \Delta t) = \varphi(t^0) + \varphi'(t^0)\Delta t + o(\Delta t).$$

ullet Так как мы хотим $t^0 + \Delta t pprox t^*$, то

$$\varphi(t^0 + \Delta t) \approx \varphi(t^*) = 0 \implies \varphi(t^0) + \varphi'(t^0) \Delta t \approx 0.$$



Задача поиска нуля: метод Ньютона

• Из $\varphi(t^0) + \varphi'(t^0) \Delta t \approx 0$ получаем:

$$\Delta t pprox rac{arphi(t^0)}{arphi'(t^0)}.$$

Задача поиска нуля: метод Ньютона

ullet Из $arphi(t^0)+arphi'(t^0)\Delta tpprox 0$ получаем:

$$\Delta t pprox rac{arphi(t^0)}{arphi'(t^0)}.$$

• Значит получаем новую точку $t^1 = t^0 + \Delta t$. Откуда получается итеративный метод:

$$t^{k+1} = t^k - rac{arphi(t^k)}{arphi'(t^k)}$$

• Этот метод называется методом Ньютона. Его предложил во второй половине 17го века тот самый Ньютон.



• **Bonpoc**: какие есть вопросы к интуиции получения итерации метода Ньютона?

• **Bonpoc**: какие есть вопросы к интуиции получения итерации метода Ньютона? Важно, что t^0 из «хорошей окрестности» t^* .

• Вопрос: какие есть вопросы к интуиции получения итерации метода Ньютона? Важно, что t^0 из «хорошей окрестности» t^* .

Рассмотрим
$$\varphi(t) = \frac{t}{\sqrt{1+t^2}}.$$
 Вопрос: какое решение? $\varphi(\phi) = 0$

$$\varphi(0) = 0$$

• **Bonpoc**: какие есть вопросы к интуиции получения итерации метода Ньютона? Важно, что t^0 из «хорошей окрестности» t^* .

• Рассмотрим

$$\varphi(t) = \frac{t}{\sqrt{1+t^2}}.$$

Вопрос: какое решение? $t^* = 0$.

- **Bonpoc**: какие есть вопросы к интуиции получения итерации метода Ньютона? Важно, что t^0 из «хорошей окрестности» t^* .
- Рассмотрим

$$\varphi(t) = \frac{t}{\sqrt{1+t^2}}.$$

Вопрос: какое решение? $t^* = 0$.

ullet Производная: $arphi'(t) = rac{1}{(1+t^2)^{3/2}}.$

• **Bonpoc**: какие есть вопросы к интуиции получения итерации метода Ньютона? Важно, что t^0 из «хорошей окрестности» t^* .

• Рассмотрим

$$\varphi(t) = \frac{t}{\sqrt{1+t^2}}.$$

Вопрос: какое решение? $t^* = 0$.

$$t^{k+1} = t^k - \frac{\varphi(t^k)}{\varphi'(t^k)} = -(t^k)^3.$$

1)
$$|t^{\circ}| < 1$$

2) $|t^{\circ}| = 1$
3) $|t^{\circ}| > 1$

- **Bonpoc**: какие есть вопросы к интуиции получения итерации метода Ньютона? Важно, что t^0 из «хорошей окрестности» t^* .
- Рассмотрим

$$\varphi(t) = \frac{t}{\sqrt{1+t^2}}.$$

Вопрос: какое решение? $t^* = 0$.

ullet Производная: $arphi'(t) = rac{1}{(1+t^2)^{3/2}}.$ Откуда итерация метода Ньютона

$$t^{k+1} = t^k - \frac{\varphi(t^k)}{\varphi'(t^k)} = -(t^k)^3.$$

• **Вопрос:** что можем сказать о сходимости к решению?

- **Bonpoc**: какие есть вопросы к интуиции получения итерации метода Ньютона? Важно, что t^0 из «хорошей окрестности» t^* .
- Рассмотрим

$$\varphi(t) = \frac{t}{\sqrt{1+t^2}}.$$

Вопрос: какое решение? $t^* = 0$.

$$t^{k+1} = t^k - \frac{\varphi(t^k)}{\varphi'(t^k)} = -(t^k)^3.$$

- Вопрос: что можем сказать о сходимости к решению?
 - ullet $|t^0| < 1$ есть сходимость

- **Bonpoc**: какие есть вопросы к интуиции получения итерации метода Ньютона? Важно, что t^0 из «хорошей окрестности» t^* .
- Рассмотрим

$$\varphi(t) = \frac{t}{\sqrt{1+t^2}}.$$

Вопрос: какое решение? $t^* = 0$.

$$t^{k+1} = t^k - \frac{\varphi(t^k)}{\varphi'(t^k)} = -(t^k)^3.$$

- Вопрос: что можем сказать о сходимости к решению?
 - ullet $|t^0| < 1$ есть сходимость
 - ullet $|t^0|=1$ колеблемся в точка -1 и 1

- **Bonpoc**: какие есть вопросы к интуиции получения итерации метода Ньютона? Важно, что t^0 из «хорошей окрестности» t^* .
- Рассмотрим

$$\varphi(t) = \frac{t}{\sqrt{1+t^2}}.$$

Вопрос: какое решение? $t^* = 0$.

$$t^{k+1} = t^k - \frac{\varphi(t^k)}{\varphi'(t^k)} = -(t^k)^3.$$

- Вопрос: что можем сказать о сходимости к решению?
 - ullet $|t^0| < 1$ есть сходимость
 - ullet $|t^0|=1$ колеблемся в точка -1 и 1
 - ullet $|t^0| > 1$ расходимся

- **Bonpoc**: какие есть вопросы к интуиции получения итерации метода Ньютона? Важно, что t^0 из «хорошей окрестности» t^* .
- Рассмотрим $\varphi(t) = \frac{t}{\sqrt{1+t^2}}.$

Вопрос: какое решение? $t^* = 0$.

$$t^{k+1} = t^k - \frac{\varphi(t^k)}{\varphi'(t^k)} = -(t^k)^3.$$

- Вопрос: что можем сказать о сходимости к решению?
 - ullet $|t^0| < 1$ есть сходимость
 - ullet $|t^0|=1$ колеблемся в точка -1 и 1
 - ullet $|t^0| > 1$ расходимся
- Ключевая особенность метода Ньютона локальная сходимость (только в окрестности решения).



Метод Ньютона: оптимизация

• Рассмотрим задачу безусловную задачу оптимизации с выпуклой дважды непрерывно дифференцируемой целевой функцией:

$$\underbrace{\min_{x \in \mathbb{R}^d} f(x)}.$$

• Вопрос: для такой задачи мы тоже ищем 0, но чего?

$$\varphi(t^*) = 0 \qquad \longrightarrow \qquad \nabla f(x^*) = 0$$

$$\chi^{(c+1)} = \chi^{k} - \left(\nabla^2 f(x^{(c)})^{-1} \nabla f(x^{(c)})\right)$$

Метод Ньютона: оптимизация

• Рассмотрим задачу безусловную задачу оптимизации с выпуклой дважды непрерывно дифференцируемой целевой функцией:

$$\min_{x \in \mathbb{R}^d} f(x)$$
.

• Вопрос: для такой задачи мы тоже ищем 0, но чего? $\nabla f(x^*) = 0$.

Метод Ньютона: оптимизация

• Рассмотрим задачу безусловную задачу оптимизации с выпуклой дважды непрерывно дифференцируемой целевой функцией:

$$\min_{x \in \mathbb{R}^d} f(x)$$
.

• **Bonpoc**: для такой задачи мы тоже ищем 0, но чего? $\nabla f(x^*) = 0$. Откуда метод Ньютона для задачи оптимизации

Алгоритм 3 Метод Ньютона

Вход: стартовая точка $x^0 \in \mathbb{R}^d$, количество итераций K

- 1: **for** k = 0, 1, ..., K 1 **do**
- 2: Вычислить $\nabla f(x^k)$, $\nabla^2 f(x^k)$
- 3: $x^{k+1} = x^k (\nabla^2 f(x^k))^{-1} \nabla f(x^k)$
- 4: end for

Выход: x^K



 Градиентный спуск работает с линейной аппроксимацией в текущей точке, метод Ньютона — с квадратичной:

$$f(x) \approx f(x^{k}) + \langle \nabla f(x^{k}), x - x^{k} \rangle + \frac{1}{2} \langle x - x^{k}, \nabla^{2} f(x^{k})(x - x^{k}) \rangle.$$

$$\nabla f(x^{k}) + \nabla^{2} f(x^{k}) (x - x^{k}) = O$$

$$\chi^{(c+1)} = \chi^{(c)} - (\nabla^{2} f(x^{(c)})^{-1} \nabla f(x^{(c)})$$

$$\chi^{(c+1)} = \chi^{(c)} - (\nabla^{2} f(x^{(c)})^{-1} \nabla f(x^{(c)})$$

• Градиентный спуск работает с линейной аппроксимацией в текущей точке, метод Ньютона — с квадратичной:

$$f(x) \approx f(x^k) + \langle \nabla f(x^k), x - x^k \rangle + \frac{1}{2} \langle x - x^k, \nabla^2 f(x^k)(x - x^k) \rangle.$$

Минимизируем квадратичную аппроксимацию по x:

 Градиентный спуск работает с линейной аппроксимацией в текущей точке, метод Ньютона — с квадратичной:

$$f(x) \approx f(x^k) + \langle \nabla f(x^k), x - x^k \rangle + \frac{1}{2} \langle x - x^k, \nabla^2 f(x^k)(x - x^k) \rangle.$$

Минимизируем квадратичную аппроксимацию по х:

$$\nabla f(x^k) = \nabla^2 f(x^k)(x - x^k) = 0.$$

 Градиентный спуск работает с линейной аппроксимацией в текущей точке, метод Ньютона — с квадратичной:

$$f(x) \approx f(x^k) + \langle \nabla f(x^k), x - x^k \rangle + \frac{1}{2} \langle x - x^k, \nabla^2 f(x^k)(x - x^k) \rangle.$$

Минимизируем квадратичную аппроксимацию по х:

 $\nabla f(x^k) = \nabla^2 f(x^k)(x - x^k) = 0$. Откуда получаем следующую точку метода:

$$x^{k+1} = x^k - \left(\nabla^2 f(x^k)\right)^{-1} \nabla f(x^k).$$

 Градиентный спуск работает с линейной аппроксимацией в текущей точке, метод Ньютона — с квадратичной:

$$f(x) \approx f(x^k) + \langle \nabla f(x^k), x - x^k \rangle + \frac{1}{2} \langle x - x^k, \nabla^2 f(x^k)(x - x^k) \rangle.$$

Минимизируем квадратичную аппроксимацию по х:

 $abla f(x^k) =
abla^2 f(x^k)(x - x^k) = 0$. Откуда получаем следующую точку метода:

 $x^{k+1} = x^k - \left(\nabla^2 f(x^k)\right)^{-1} \nabla f(x^k).$

• Метод Ньютона использует оракул второго порядка: требует вычисление гессиана.

 Градиентный спуск работает с линейной аппроксимацией в текущей точке, метод Ньютона — с квадратичной:

$$f(x) \approx f(x^k) + \langle \nabla f(x^k), x - x^k \rangle + \frac{1}{2} \langle x - x^k, \nabla^2 f(x^k)(x - x^k) \rangle.$$

Минимизируем квадратичную аппроксимацию по х:

 $\nabla f(x^k) = \nabla^2 f(x^k)(x - x^k) = 0$. Откуда получаем следующую точку метода:

 $x^{k+1} = x^k - \left(\nabla^2 f(x^k)\right)^{-1} \nabla f(x^k).$

- Метод Ньютона использует оракул второго порядка: требует вычисление гессиана.
- Стоимость итерации значительно возрастает (по сравнению с градиентным спуском) не только из-за гессиана, но и его обращения.

 Градиентный спуск работает с линейной аппроксимацией в текущей точке, метод Ньютона — с квадратичной:

$$f(x) \approx f(x^k) + \langle \nabla f(x^k), x - x^k \rangle + \frac{1}{2} \langle x - x^k, \nabla^2 f(x^k)(x - x^k) \rangle.$$

Минимизируем квадратичную аппроксимацию по х:

 $abla f(x^k) =
abla^2 f(x^k)(x - x^k) = 0$. Откуда получаем следующую точку метода:

 $x^{k+1} = x^k - \left(\nabla^2 f(x^k)\right)^{-1} \nabla f(x^k). \quad \text{if } X \leftarrow X \leftarrow D \times C$

- Метод Ньютона использует оракул второго порядка: требует вычисление гессиана. $\chi^{lef1} = \chi^{lc} A^{-1}(A\chi^{lc} b) = A^{-1}b$
- Стоимость итерации значительно возрастает (по сравнению с градиентным спуском) не только из-за гессиана, но и его обращения. Вопрос: за сколько итераций метод Ньютона сойдется для квадратичной задачи с положительно определенной матрицей?

 Градиентный спуск работает с линейной аппроксимацией в текущей точке, метод Ньютона — с квадратичной:

$$f(x) \approx f(x^k) + \langle \nabla f(x^k), x - x^k \rangle + \frac{1}{2} \langle x - x^k, \nabla^2 f(x^k)(x - x^k) \rangle.$$

Минимизируем квадратичную аппроксимацию по х:

 $abla f(x^k) =
abla^2 f(x^k)(x - x^k) = 0$. Откуда получаем следующую точку метода:

 $x^{k+1} = x^k - \left(\nabla^2 f(x^k)\right)^{-1} \nabla f(x^k).$

- Метод Ньютона использует оракул второго порядка: требует вычисление гессиана.
- Стоимость итерации значительно возрастает (по сравнению с градиентным спуском) не только из-за гессиана, но и его обращения. Вопрос: за сколько итераций метод Ньютона сойдется для квадратичной задачи с положительно определенной матрицей? за 1 (но дорогую).

 То, что для квадратичной задачи метод Ньютона сходится за 1 итерацию, наталкивает на мысль о том, что при всех своих минусах (локальная сходимость, дороговизна итерации) ключевым плюсом является скорость сходимости.

- То, что для квадратичной задачи метод Ньютона сходится за 1 итерацию, наталкивает на мысль о том, что при всех своих минусах (локальная сходимость, дороговизна итерации) ключевым плюсом является скорость сходимости.
- Пусть целевая функция в задаче безусловной минимизации является дважды непрерывно дифференцируемой, μ -сильно выпуклой и имеет M-Липшицев гессиан, т.е. для любых $x, y \in \mathbb{R}^d$ справедливо:

$$\nabla^2 f(x) \succeq \mu I$$
, $\|\nabla^2 f(x) - \nabla^2 f(y)\|_2 \leq M\|x - y\|_2$.

В случае матрицы $\|\cdot\|_2$ — спектральная норма (согласованная норма с евклидовой для векторов). $\|x\| = 1$

$$\| \triangle_{\mathcal{A}}(X) \|^{5} \geq \mathbb{I}$$

• Доказываем сходимость.

 Доказываем сходимость. Будем изучать, как меняется расстояние до решения:

$$x^{k+1} - x^* = x^k - \left(\nabla^2 f(x^k)\right)^{-1} \nabla f(x^k) - x^*.$$

$$= \int \nabla^2 f(x^k) - \nabla f(x^k) = \int \nabla^2 f(x^k + \tau(x^k - x^k)) \left(x^{k} - x^k\right) d\tau$$

$$= \int \int \nabla^2 f(x^k) - \left(x^k - x^k\right) d\tau d\tau$$

$$= \int \int \nabla^2 f(x^k) - \left(x^k - x^k\right) d\tau d\tau$$

$$= \int \int \int \nabla^2 f(x^k + \tau(x^k - x^k)) d\tau d\tau$$

$$= \int \int \int \int \nabla^2 f(x^k - x^k) d\tau d\tau$$

$$= \int \int \int \int \int \nabla^2 f(x^k - x^k) d\tau d\tau$$

• Доказываем сходимость. Будем изучать, как меняется расстояние до решения:

$$x^{k+1} - x^* = x^k - \left(\nabla^2 f(x^k)\right)^{-1} \nabla f(x^k) - x^*.$$

 Снова вспомним формулу Ньютона-Лейбница для интеграла вдоль кривой:

$$\nabla f(x^{k}) - \nabla f(x^{*}) = \int_{0}^{1} \nabla^{2} f(x^{*} + \tau(x^{k} - x^{*}))(x^{k} - x^{*}) d\tau$$

• Доказываем сходимость. Будем изучать, как меняется расстояние до решения:

$$x^{k+1} - x^* = x^k - \left(\nabla^2 f(x^k)\right)^{-1} \nabla f(x^k) - x^*.$$

 Снова вспомним формулу Ньютона-Лейбница для интеграла вдоль кривой:

$$\nabla f(x^k) - \nabla f(x^*) = \int_0^1 \nabla^2 f(x^* + \tau(x^k - x^*))(x^k - x^*) d\tau$$

Зная, что $\nabla f(x^*) = 0$, получим

$$x^{k+1}-x^*=x^k-x^*-\left(\nabla^2 f(x^k)\right)^{-1}\int_0^1 \nabla^2 f(x^*+\tau(x^k-x^*))(x^k-x^*)d\tau.$$



• Продолжаем и используем «умную единицу»:

$$x^{k+1} - x^* = x^k - x^* - \left(\nabla^2 f(x^k)\right)^{-1} \int_0^1 \nabla^2 f(x^* + \tau(x^k - x^*))(x^k - x^*)$$

$$= \left(\nabla^2 f(x^k)\right)^{-1} \nabla^2 f(x^k)(x^k - x^*)$$

$$- \left(\nabla^2 f(x^k)\right)^{-1} \int_0^1 \nabla^2 f(x^* + \tau(x^k - x^*))(x^k - x^*) d\tau.$$

• Продолжаем и используем «умную единицу»:

$$x^{k+1} - x^* = x^k - x^* - \left(\nabla^2 f(x^k)\right)^{-1} \int_0^1 \nabla^2 f(x^* + \tau(x^k - x^*))(x^k - x^*)$$

$$= \left(\nabla^2 f(x^k)\right)^{-1} \nabla^2 f(x^k)(x^k - x^*)$$

$$- \left(\nabla^2 f(x^k)\right)^{-1} \int_0^1 \nabla^2 f(x^* + \tau(x^k - x^*))(x^k - x^*) d\tau.$$

• Заметим, что $x^k - x^*$ можно вынести за пределы интеграла:

$$x^{k+1} - x^* = \left(\nabla^2 f(x^k)\right)^{-1} \nabla^2 f(x^k) (x^k - x^*)$$
$$- \left(\nabla^2 f(x^k)\right)^{-1} \left(\int_0^1 \nabla^2 f(x^* + \tau(x^k - x^*)) d\tau\right) (x^k - x^*)$$

ullet Введем обозначение $G_k =
abla^2 f(x^k) - \int_0^1
abla^2 f(x^* + au(x^k - x^*)) d au$:

$$x^{k+1} - x^* = \left(\nabla^2 f(x^k)\right)^{-1} G_k(x^k - x^*).$$

$$\|x^{k+1} - x^*\|_2 = \left\|\left(\nabla^2 f(x^k)\right)^{-1} G_k\left(x^k - x^*\right)\right\|_2$$

$$\leq \left\|\left(\nabla^2 f(x^k)\right)^{-1} G_k\left\|_2 \left\|x^k - x^*\right\|_2$$

$$\leq \left\|\left(\nabla^2 f(x^k)\right)^{-1}\right\|_2 \left\|G_k\right\|_2 \left\|x^k - x^*\right\|_2$$

$$\leq \left\|\left(\nabla^2 f(x^k)\right)^{-1}\right\|_2 \left\|G_k\right\|_2 \left\|x^k - x^*\right\|_2$$

$$\leq \left\|\left(G_k\right)^{-1}\right\|_2 \left\|G_k\right\|_2 \left\|x^k - x^*\right\|_2$$

• Введем обозначение $G_k = \nabla^2 f(x^k) - \int_0^1 \nabla^2 f(x^* + \tau(x^k - x^*)) d\tau$:

$$x^{k+1} - x^* = \left(\nabla^2 f(x^k)\right)^{-1} G_k(x^k - x^*).$$

• Перейдем к оценке нормы расстояния:

$$\|x^{k+1} - x^*\|_2 = \left\| \left(\nabla^2 f(x^k) \right)^{-1} G_k(x^k - x^*) \right\|_2$$

• Введем обозначение $G_k = \nabla^2 f(x^k) - \int_0^1 \nabla^2 f(x^* + \tau(x^k - x^*)) d\tau$:

$$x^{k+1} - x^* = \left(\nabla^2 f(x^k)\right)^{-1} G_k(x^k - x^*).$$

• Перейдем к оценке нормы расстояния:

$$||x^{k+1} - x^*||_2 = \left\| \left(\nabla^2 f(x^k) \right)^{-1} G_k(x^k - x^*) \right\|_2$$

• Пользуемся, что спектральная норма матрицы согласована с евклидовой вектора:

$$||x^{k+1} - x^*||_2 \le ||\left(\nabla^2 f(x^k)\right)^{-1} G_k||_2 ||x^k - x^*||_2$$

$$\le ||\left(\nabla^2 f(x^k)\right)^{-1}||_2 ||G_k||_2 ||x^k - x^*||_2.$$



• С предыдущего слайда:

$$||x^{k+1}-x^*||_2 \le ||(\nabla^2 f(x^k))^{-1}||_2 ||G_k||_2 ||x^k-x^*||_2.$$

• С предыдущего слайда:

$$||x^{k+1}-x^*||_2 \le ||(\nabla^2 f(x^k))^{-1}||_2 ||G_k||_2 ||x^k-x^*||_2.$$

• Вопрос: как оценить $\| (\nabla^2 f(x^k))^{-1} \|_2$?

• С предыдущего слайда:

$$||x^{k+1}-x^*||_2 \le ||(\nabla^2 f(x^k))^{-1}||_2 ||G_k||_2 ||x^k-x^*||_2.$$

• Вопрос: как оценить $\left\| \left(\nabla^2 f(x^k) \right)^{-1} \right\|_2$? Мы знаем, что $\nabla^2 f(x) \succeq \mu I$, а значит $\frac{1}{\mu} I \succeq \left(\nabla^2 f(x^k) \right)^{-1}$,

• С предыдущего слайда:

$$||x^{k+1}-x^*||_2 \le ||(\nabla^2 f(x^k))^{-1}||_2 ||G_k||_2 ||x^k-x^*||_2.$$

• Вопрос: как оценить $\left\| \left(\nabla^2 f(x^k) \right)^{-1} \right\|_2$? Мы знаем, что $\nabla^2 f(x) \succeq \mu I$, а значит $\frac{1}{\mu} I \succeq \left(\nabla^2 f(x^k) \right)^{-1}$, откуда $\left\| \left(\nabla^2 f(x^k) \right)^{-1} \right\|_2 \leq \frac{1}{\mu}$ и

$$||x^{k+1} - x^*||_2 \le \frac{1}{\mu} ||G_k||_2 ||x^k - x^*||_2.$$

• С предыдущего слайда:

$$||x^{k+1}-x^*||_2 \le ||(\nabla^2 f(x^k))^{-1}||_2 ||G_k||_2 ||x^k-x^*||_2.$$

• Вопрос: как оценить $\left\| \left(\nabla^2 f(x^k) \right)^{-1} \right\|_2$? Мы знаем, что $\nabla^2 f(x) \succeq \mu I$, а значит $\frac{1}{\mu} I \succeq \left(\nabla^2 f(x^k) \right)^{-1}$, откуда $\left\| \left(\nabla^2 f(x^k) \right)^{-1} \right\|_2 \leq \frac{1}{\mu}$ и

$$||x^{k+1} - x^*||_2 \le \frac{1}{\mu} ||G_k||_2 ||x^k - x^*||_2.$$

• Осталось оценить $||G_k||_2$.



$$||G_{k}||_{2} = ||\nabla^{2}f(x^{k}) - \int_{0}^{2}\nabla^{2}f(x^{k} + \tau(x^{k} - x^{k})) d\tau||_{2}$$

$$= ||\int_{0}^{2}(\nabla^{2}f(x^{k}) - \nabla^{2}f(x^{k} + \tau(x^{k} - x^{k}))) d\tau||_{2}$$

$$\leq \int_{0}^{2}||\nabla^{2}f(x^{k}) - \nabla^{2}f(x^{k} + \tau(x^{k} - x^{k}))||_{2} d\tau$$

$$\leq \int_{0}^{2}||\nabla^{2}f(x^{k}) - \nabla^{2}f(x^{k} + \tau(x^{k} - x^{k}))||_{2} d\tau$$

$$\leq \int_{0}^{2}||\nabla^{2}f(x^{k}) - \nabla^{2}f(x^{k} + \tau(x^{k} - x^{k}))||_{2} d\tau$$

$$\leq \int_{0}^{2}||\nabla^{2}f(x^{k}) - \nabla^{2}f(x^{k} + \tau(x^{k} - x^{k}))||_{2} d\tau$$

$$= \int_{0}^{2}||\nabla^{2}f(x^{k}) - \nabla^{2}f(x^{k} + \tau(x^{k} - x^{k}))||_{2} d\tau$$

$$= \int_{0}^{2}||\nabla^{2}f(x^{k}) - \nabla^{2}f(x^{k} + \tau(x^{k} - x^{k}))||_{2} d\tau$$

$$= \int_{0}^{2}||\nabla^{2}f(x^{k}) - \nabla^{2}f(x^{k} + \tau(x^{k} - x^{k}))||_{2} d\tau$$

$$= \int_{0}^{2}||\nabla^{2}f(x^{k}) - \nabla^{2}f(x^{k} + \tau(x^{k} - x^{k}))||_{2} d\tau$$

$$= \int_{0}^{2}||\nabla^{2}f(x^{k}) - \nabla^{2}f(x^{k} + \tau(x^{k} - x^{k}))||_{2} d\tau$$

$$= \int_{0}^{2}||\nabla^{2}f(x^{k}) - \nabla^{2}f(x^{k} + \tau(x^{k} - x^{k}))||_{2} d\tau$$

$$= \int_{0}^{2}||\nabla^{2}f(x^{k}) - \nabla^{2}f(x^{k} + \tau(x^{k} - x^{k}))||_{2} d\tau$$

$$= \int_{0}^{2}||\nabla^{2}f(x^{k}) - \nabla^{2}f(x^{k} + \tau(x^{k} - x^{k}))||_{2} d\tau$$

$$= \int_{0}^{2}||\nabla^{2}f(x^{k}) - \nabla^{2}f(x^{k} + \tau(x^{k} - x^{k}))||_{2} d\tau$$

$$= \int_{0}^{2}||\nabla^{2}f(x^{k}) - \nabla^{2}f(x^{k} + \tau(x^{k} - x^{k}))||_{2} d\tau$$

$$= \int_{0}^{2}||\nabla^{2}f(x^{k}) - \nabla^{2}f(x^{k} + \tau(x^{k} - x^{k}))||_{2} d\tau$$

$$= \int_{0}^{2}||\nabla^{2}f(x^{k} - x^{k})||_{2} d\tau$$

• Оцениваем $||G_k||_2$:

$$||G_{k}||_{2} = ||\nabla^{2}f(x^{k}) - \int_{0}^{1} \nabla^{2}f(x^{*} + \tau(x^{k} - x^{*}))d\tau||_{2}$$

$$= ||\int_{0}^{1} (\nabla^{2}f(x^{k}) - \nabla^{2}f(x^{*} + \tau(x^{k} - x^{*}))) d\tau||_{2}$$

$$\leq \int_{0}^{1} ||\nabla^{2}f(x^{k}) - \nabla^{2}f(x^{*} + \tau(x^{k} - x^{*}))||_{2} d\tau$$

$$\leq \int_{0}^{1} M(1 - \tau)||x^{k} - x^{*}||_{2} d\tau$$

$$= M||x^{k} - x^{*}||_{2} \int_{0}^{1} (1 - \tau)d\tau = \frac{M}{2}||x^{k} - x^{*}||_{2}.$$

• Подставляем оценку на $\|G_k\|_2$:

$$||x^{k+1} - x^*||_2 \le \frac{M}{2\mu} ||x^k - x^*||_2^2.$$

• Подставляем оценку на $\|G_k\|_2$:

$$||x^{k+1} - x^*||_2^0 \le \frac{M}{2\mu} ||x^k - x^*||_2^2$$

Теорема об оценке сходимости метода Ньютона для μ -сильно выпуклых функций с M-Липшецевым гессианом

Пусть задача безусловной оптимизации с μ -сильно выпуклой целевой функцией f с M-Липшецевыми гессианом решается методом Ньютона. Тогда справедлива следующая оценка сходимости за 1 итерацию

$$||x^{k+1} - x^*||_2 \le \frac{M}{2\mu} ||x^k - x^*||_2^2.$$



ullet Подставляем оценку на $\|G_k\|_2$:

$$||x^{k+1} - x^*||_2 \le \frac{M}{2\mu} ||x^k - x^*||_2^2.$$

Теорема об оценке сходимости метода Ньютона для μ -сильно выпуклых функций с M-Липшецевым гессианом

Пусть задача безусловной оптимизации с μ -сильно выпуклой целевой функцией f с M-Липшецевыми гессианом решается методом Ньютона. Тогда справедлива следующая оценка сходимости за 1 итерацию

$$||x^{k+1} - x^*||_2 \le \frac{M}{2\mu} ||x^k - x^*||_2^2.$$

Мы уже знаем, что такого рода оценки дают квадратичную скорость сходимости.



• Сходимость, как и в случае первородного метода Ньютона, является локальной.

• Сходимость, как и в случае первородного метода Ньютона, является локальной. А именно, что гарантировать $\|x^1-x^*\|_2<\|x^0-x^*\|_2$ нужно предположить, что

$$||x^0 - x^*||_2 < \frac{2\mu}{M}.$$

• Сходимость, как и в случае первородного метода Ньютона, является локальной. А именно, что гарантировать $\|x^1-x^*\|_2<\|x^0-x^*\|_2$ нужно предположить, что

$$||x^0 - x^*||_2 < \frac{2\mu}{M}.$$

• Поймем насколько быстро сходится метод. Пусть M=2, $\mu=1$, а $\|x^0-x^*\|_2=\frac{1}{2}$.

• Сходимость, как и в случае первородного метода Ньютона, является локальной. А именно, что гарантировать $\|x^1-x^*\|_2<\|x^0-x^*\|_2$ нужно предположить, что

$$||x^0 - x^*||_2 < \frac{2\mu}{M}.$$

• Поймем насколько быстро сходится метод. Пусть M=2, $\mu=1$, а $\|x^0-x^*\|_2=\frac{1}{2}$. Тогда мы можем гарантировать, что $\|x^1-x^*\|_2\leq \frac{1}{2^2}$,

• Сходимость, как и в случае первородного метода Ньютона, является локальной. А именно, что гарантировать $\|x^1-x^*\|_2<\|x^0-x^*\|_2$ нужно предположить, что

$$||x^0 - x^*||_2 < \frac{2\mu}{M}.$$

• Поймем насколько быстро сходится метод. Пусть M=2, $\mu=1$, а $\|x^0-x^*\|_2=\frac{1}{2}$. Тогда мы можем гарантировать, что $\|x^1-x^*\|_2\leq \frac{1}{2^2}$, $\|x^2-x^*\|_2\leq \frac{1}{(2^2)^2}$ и так далее.



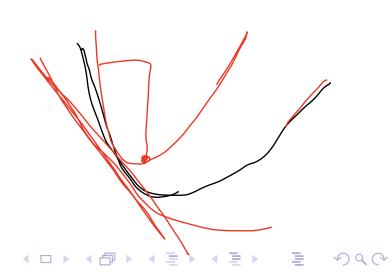
• Пытаемся решить проблему локальной сходимости. Действуем по аналогии с градиентным спуском. Вопрос: идеи?

• Пытаемся решить проблему локальной сходимости. Действуем по аналогии с градиентным спуском. **Вопрос**: идеи?

• Идея первая – шаг: $x^{k+1} = x^k - \gamma_k \left(\nabla^2 f(x^k)\right)^{-1} \nabla f(x^k).$

Такой метод называется демпфированный метод Ньютона.

$$X^{(c+1)} = X^k + X^k P^k$$



- Пытаемся решить проблему локальной сходимости. Действуем по аналогии с градиентным спуском. Вопрос: идеи?
- Идея первая шаг:

$$x^{k+1} = x^k - \gamma_k \left(\nabla^2 f(x^k) \right)^{-1} \nabla f(x^k).$$

Такой метод называется демпфированный метод Ньютона.

Вопрос: как выбирать шаг?



- Пытаемся решить проблему локальной сходимости. Действуем по аналогии с градиентным спуском. Вопрос: идеи?
- Идея первая шаг:

$$x^{k+1} = x^k - \gamma_k \left(\nabla^2 f(x^k) \right)^{-1} \nabla f(x^k).$$

Такой метод называется демпфированный метод Ньютона. **Вопрос:** как выбирать шаг? Много разных способов, например, на прошлой лекции обсуждали линейный поиск: arg $\min_{\gamma}(x^k + \gamma p_k)$, где $p_k = -\left(\nabla^2 f(x^k)\right)^{-1} \nabla f(x^k)$.

• Идея вторая – «оценки сверху». В основе анализа градиентного спуска лежала оптимизация «оценки сверху» на функцию:

$$x^{k+1} = \arg\min_{x \in \mathbb{R}^d} \left(f(x^k) + \langle \nabla f(x^k), x - x^k \rangle + \frac{L}{2} \|x - x^k\|_2^2 \right).$$

• Идея вторая – «оценки сверху». В основе анализа градиентного спуска лежала оптимизация «оценки сверху» на функцию:

$$x^{k+1} = \arg\min_{x \in \mathbb{R}^d} \left(f(x^k) + \langle \nabla f(x^k), x - x^k \rangle + \frac{L}{2} \|x - x^k\|_2^2 \right).$$

Вопрос: чему равно x^{k+1} ?

• Идея вторая – «оценки сверху». В основе анализа градиентного спуска лежала оптимизация «оценки сверху» на функцию:

$$x^{k+1} = \arg\min_{x \in \mathbb{R}^d} \left(f(x^k) + \langle \nabla f(x^k), x - x^k \rangle + \frac{L}{2} \|x - x^k\|_2^2 \right).$$

Вопрос: чему равно x^{k+1} ? $x^{k+1} = x^k - \frac{1}{L}\nabla f(x^k)$.

• Идея вторая – «оценки сверху». В основе анализа градиентного спуска лежала оптимизация «оценки сверху» на функцию:

$$x^{k+1} = \arg\min_{x \in \mathbb{R}^d} \left(f(x^k) + \langle \nabla f(x^k), x - x^k \rangle + \frac{L}{2} \|x - x^k\|_2^2 \right).$$

Вопрос: чему равно x^{k+1} ? $x^{k+1} = x^k - \frac{1}{L}\nabla f(x^k)$. Запишем, похожее для аппроксимации 2-го порядка:

$$x^{k+1} = \arg\min_{x \in \mathbb{R}^d} \left(f(x^k) + \langle \nabla f(x^k), x - x^k \rangle + \frac{1}{2} \langle x - x^k, \nabla^2 f(x^k)(x - x^k) \rangle + \frac{M}{6} ||x - x^k||_2^3 \right).$$

Здесь M — константа Липшица гессиана. Такой метод называется кубический метод Ньютона.



• Запишем метод Ньютона следующим образом:

$$x^{k+1} = x^k - \gamma_k H_k \nabla f(x^k).$$

$$(\nabla^k \int (x^k)^{-1}$$

• Запишем метод Ньютона следующим образом:

$$x^{k+1} = x^k - \gamma_k H_k \nabla f(x^k).$$

В случае метода Ньютона вместо H_k стоит $(\nabla^2 f(x^k))^{-1}$.

• Хочется заменить $\left(\nabla^2 f(x^k)\right)^{-1}$ на что-то более дешевое с точки зрения вычислений.

• Запишем метод Ньютона следующим образом:

$$x^{k+1} = x^k - \gamma_k H_k \nabla f(x^k).$$

В случае метода Ньютона вместо H_k стоит $(\nabla^2 f(x^k))^{-1}$.

- Хочется заменить $\left(\nabla^2 f(x^k)\right)^{-1}$ на что-то более дешевое с точки зрения вычислений.
- Идея выудить какие-то свойства присущие гессиану.

• Запишем метод Ньютона следующим образом:

$$x^{k+1} = x^k - \gamma_k H_k \nabla f(x^k).$$

В случае метода Ньютона вместо H_k стоит $(\nabla^2 f(x^k))^{-1}$.

- Хочется заменить $\left(\nabla^2 f(x^k) \right)^{-1}$ на что-то более дешевое с точки зрения вычислений.
- Идея выудить какие-то свойства присущие гессиану.
- Связь градиента и гессиана:

$$abla f(x^k) =
abla f(x^{k+1}) +
abla^2 f(x^{k+1})(x^k - x^{k+1}) + o(||x^{k+1} - x^k||_2)$$
или $abla f(x^k) -
abla f(x^{k+1}) \approx
abla^2 f(x^{k+1})(x^k - x^{k+1}).$ Откуда $abla^{k+1} - x^k \approx (
abla^2 f(x^{k+1}))^{-1} (
abla f(x^{k+1}) -
abla f(x^k) + o(||x^{k+1} - x^k||_2).$

• Запишем метод Ньютона следующим образом:

$$x^{k+1} = x^k - \gamma_k H_k \nabla f(x^k).$$

В случае метода Ньютона вместо H_k стоит $(\nabla^2 f(x^k))^{-1}$.

- Хочется заменить $\left(\nabla^2 f(x^k)\right)^{-1}$ на что-то более дешевое с точки зрения вычислений.
- Идея выудить какие-то свойства присущие гессиану.
- Связь градиента и гессиана:

$$\nabla f(x^k) = \nabla f(x^{k+1}) + \nabla^2 f(x^{k+1})(x^k - x^{k+1}) + o(\|x^{k+1} - x^k\|_2)$$

или $\nabla f(x^k) - \nabla f(x^{k+1}) \approx \nabla^2 f(x^{k+1})(x^k - x^{k+1})$. Откуда $x^{k+1} - x^k \approx (\nabla^2 f(x^{k+1}))^{-1}(\nabla f(x^{k+1}) - \nabla f(x^k))$. Заменим на $(\nabla^2 f(x^{k+1}))^{-1}$ на H_{k+1} , введем обозначения $\underline{s^k = x^{k+1} - x^k}$ и $y^k = \nabla f(x^{k+1}) - \nabla f(x^k)$:

$$s^k = H_{k+1}y^k$$



• Квазиньютоновское уравнение:

$$s^k = H_{k+1} y^k$$

• Квазиньютоновское уравнение:

$$s^k = H_{k+1} y^k$$

ullet Еще потребуем, чтобы H_{k+1} была симметричной: $H_{k+1}^{T} = H_{k+1}$.

$$S^{k} = H_{lefn}g^{k}$$
d ypabrennin $O(d^{2})$

• Квазиньютоновское уравнение:

$$s^k = H_{k+1} y^k$$

- ullet Еще потребуем, чтобы H_{k+1} была симметричной: $H_{k+1}^T = H_{k+1}$.
- Вопрос: сколько решений имеет система уравнений $s^k = H_{k+1} y^k$ относительно H_{k+1} при условии, что $H_{k+1}^T = H_{k+1}$?

• Квазиньютоновское уравнение:

$$s^k = H_{k+1} y^k$$

- ullet Еще потребуем, чтобы H_{k+1} была симметричной: $H_{k+1}^{\mathcal{T}} = H_{k+1}$.
- Вопрос: сколько решений имеет система уравнений $s^k = H_{k+1} y^k$ относительно H_{k+1} при условии, что $H_{k+1}^T = H_{k+1}$? d уравнений, d + d(d-1)/2 уравнений. Можно урешаться.

Квазиньютоновское уравнение

• Квазиньютоновское уравнение:

$$s^k = H_{k+1} y^k$$

- ullet Еще потребуем, чтобы H_{k+1} была симметричной: $H_{k+1}^T = H_{k+1}$.
- Вопрос: сколько решений имеет система уравнений $s^k = H_{k+1} y^k$ относительно H_{k+1} при условии, что $H_{k+1}^T = H_{k+1}$? d уравнений, d + d(d-1)/2 уравнений. Можно урешаться. Нужно еще сузить правила поиска H_{k+1} .

 Идея первая – одно-ранговая (дешевая с точки зрения вычислений) добавка:

$$H_{k+1}=H_k+\mu_kq^k(q^k)^T,$$
где $\mu_k\in\mathbb{R}$ и $q^k\in\mathbb{R}^d$ нужно подобрать.

$$S^{k} = H_{k+1} y^{k} = H_{k} y^{k} + M_{k} q^{k} q^{k} y^{k} y^{k}$$

$$M_{k} (g^{k})^{T} y^{k} \qquad Q^{k} = S^{k} - H_{k} y^{k}$$

$$Q^{k} = S^{k} - H_{k} y^{k}$$

• Идея первая – одно-ранговая (дешевая с точки зрения вычислений) добавка:

$$H_{k+1} = H_k + \mu_k q^k (q^k)^T,$$

где $\mu_k \in \mathbb{R}$ и $q^k \in \mathbb{R}^d$ нужно подобрать.

• Подбираем исходя из квазиньютоновского уравнения:

$$s^{k} = H_{k+1}y^{k} = H_{k}y^{k} + \mu_{k}q^{k}(q^{k})^{T}y^{k}$$

= $H_{k}y^{k} - \mu_{k}\left((q^{k})^{T}y^{k}\right)q^{k}$

 Идея первая – одно-ранговая (дешевая с точки зрения вычислений) добавка:

$$H_{k+1} = H_k + \mu_k q^k (q^k)^T,$$

где $\mu_k \in \mathbb{R}$ и $q^k \in \mathbb{R}^d$ нужно подобрать.

• Подбираем исходя из квазиньютоновского уравнения:

$$s^{k} = H_{k+1}y^{k} = H_{k}y^{k} + \mu_{k}q^{k}(q^{k})^{T}y^{k}$$

= $H_{k}y^{k} - \mu_{k}\left((q^{k})^{T}y^{k}\right)q^{k}$

Откуда

$$\mu_k\left((q^k)^T y^k\right) q^k = s^k - H_k y^k$$



• С предыдущего слайда:

$$\mu_k\left((q^k)^T y^k\right) q^k = s^k - H_k y^k$$

• Вопрос: что можно сказать про вектор q^k ?

• С предыдущего слайда:

$$\mu_k\left((q^k)^T y^k\right) q^k = s^k - H_k y^k$$

• **Вопрос**: что можно сказать про вектор q^k ? Коллинеарен $s^k - H_k y^k$.

• С предыдущего слайда:

$$\mu_k\left((q^k)^T y^k\right) q^k = s^k - H_k y^k$$

• **Вопрос:** что можно сказать про вектор q^k ? Коллинеарен $s^k - H_k y^k$. Пусть

$$q^k = s^k - H_k y^k,$$

тогда

$$\mu_k = \frac{1}{(q^k)^T y^k}.$$

• Получаем SR1 способ подсчета матриц H:

$$H_{k+1} = H_k + \frac{(s^k - H_k y^k)(s^k - H_k y^k)^T}{(s^k - H_k y^k)^T y^k}$$





• Посмотрим на задачу поиска H_{k+1} , как на задачу поиска «близкой» к H_k матрицы с точки зрения оптимизации:

$$H_{k+1} = \arg\min_{H \in \mathbb{R}^{d \times d}} \|H - H_k\|^2$$

 $s.t. \ s^k = Hy^k$
 $H^T = H$

• Посмотрим на задачу поиска H_{k+1} , как на задачу поиска «близкой» к H_k матрицы с точки зрения оптимизации:

$$H_{k+1} = \arg\min_{H \in \mathbb{R}^{d \times d}} \|H - H_k\|^2$$

 $s.t. \ s^k = Hy^k$
 $H^T = H$

• Норма в задачи оптимизации может быть любая. В зависимости от нормы будет получаться разные квазиньютоновские методы.

• Посмотрим на задачу поиска H_{k+1} , как на задачу поиска «близкой» к H_k матрицы с точки зрения оптимизации:

$$H_{k+1} = \arg\min_{H \in \mathbb{R}^{d \times d}} \|H - H_k\|^2$$

 $s.t. \ s^k = Hy^k$
 $H^T = H$

- Норма в задачи оптимизации может быть любая. В зависимости от нормы будет получаться разные квазиньютоновские методы.
- Рассмотрим взвешенную норму Фробениуса $\|A\|_W = \|W^{1/2}AW^{1/2}\|_F$, где должно выполняться $Wy^k = s^k$. Такой выбор дает метод BFGS:

$$H_{k+1} = (I - \rho_k s^k (y^k)^T) H_k (I - \rho_k y^k (s^k)^T) + \rho_k s^k (s^k)^T$$
, где $\rho_k = \frac{1}{(y^k)^T s^k}$



• До такой формулы можно дойти по-другому.

• До такой формулы можно дойти по-другому. Рассмотрим $B_{k+1} = H_{k+1}^{-1}$. Для B квазиньютоновское уравнение выглядит как

$$B_{k+1}s^k = y^k$$

• До такой формулы можно дойти по-другому. Рассмотрим $B_{k+1} = H_{k+1}^{-1}$. Для B квазиньютоновское уравнение выглядит как

$$B_{k+1}s^k=y^k$$

• Для B_{k+1} можно написать SR1 пересчет матрицы:

$$B_{k+1} = B_k + \frac{(y^k) - (B_k s^k)(y^k) - (B_k s^k)^T}{(y^k - B_k s^k)^T s^k}$$



ullet До такой формулы можно дойти по-другому. Рассмотрим $B_{k+1} = H_{k+1}^{-1}$. Для B квазиньютоновское уравнение выглядит как

$$B_{k+1}s^k=y^k$$

• Для B_{k+1} можно написать SR1 пересчет матрицы:

$$B_{k+1} = B_k + \frac{(y^k - B_k s^k)(y^k - B_k s^k)^T}{(y^k - B_k s^k)^T s^k}$$

• Смотрим на вид B_{k+1} и делаем из нее двухранговое изменение:

$$B_{k+1} = B_k + \mu_{k,1} y^k (y^k)^T + \mu_{k,2} B_k y^k (B_k y^k)^T$$

• До такой формулы можно дойти по-другому. Рассмотрим $B_{k+1} = H_{k+1}^{-1}$. Для B квазиньютоновское уравнение выглядит как

$$B_{k+1}s^k=y^k$$

• Для B_{k+1} можно написать SR1 пересчет матрицы:

$$B_{k+1} = B_k + \frac{(y^k - B_k s^k)(y^k - B_k s^k)^T}{(y^k - B_k s^k)^T s^k}$$

ullet Смотрим на вид B_{k+1} и делаем из нее двухранговое изменение:

$$B_{k+1} = B_k + \mu_{k,1} y^k (y^k)^T + \mu_{k,2} B_k y^k (B_k y^k)^T$$

• Как и в SR1 можно подогнать $\mu_{k,1}$ и $\mu_{k,2}$:

$$B_{k+1} = B_k + \frac{y^k (y^k)^T}{(y^k)^T s^k} + \frac{B_k y^k (B_k y^k)^T}{(s^k)^T B_k s^k}$$



• До такой формулы можно дойти по-другому. Рассмотрим $B_{k+1} = H_{k+1}^{-1}$. Для B квазиньютоновское уравнение выглядит как

$$B_{k+1}s^k=y^k$$

• Для B_{k+1} можно написать SR1 пересчет матрицы:

$$B_{k+1} = B_k + \frac{(y^k - B_k s^k)(y^k - B_k s^k)^T}{(y^k - B_k s^k)^T s^k}$$

ullet Смотрим на вид B_{k+1} и делаем из нее двухранговое изменение:

$$B_{k+1} = B_k + \mu_{k,1} y^k (y^k)^T + \mu_{k,2} B_k y^k (B_k y^k)^T$$

• Как и в SR1 можно подогнать $\mu_{k,1}$ и $\mu_{k,2}$:

$$B_{k+1} = B_k + \frac{y^k (y^k)^T}{(y^k)^T s^k} + \frac{B_k y^k (B_k y^k)^T}{(s^k)^T B_k s^k}$$

• Если теперь обратить B_{k+1} (формула Шермана-Маррисона-Вудберри), то получится выражение для H_{k+1}

• Вопрос: чтобы посчитать новую матрицу нужно $O(d^2)$ операций (не учитывая подсчет градиентов). Кажется, что BFGS подсчет дороже (есть перемнножение трех матриц). Так ли это? $H_{k+1} = (I - \rho_k s^k (y^k)^T) \widehat{H}_k (I - \rho_k y^k (s^k)^T) + \rho_k s^k (s^k)^T$

$$\frac{(1 + pks)(y)}{(1 + pks)(y)} + pks(y)$$

- Вопрос: чтобы посчитать новую матрицу нужно $O(d^2)$ операций (не учитывая подсчет градиентов). Кажется, что BFGS подсчет дороже (есть перемнножение трех матриц). Так ли это? $H_{k+1} = (I \rho_k s^k (y^k)^T) H_k (I \rho_k y^k (s^k)^T) + \rho_k s^k (s^k)^T$
- Нужно раскрыть скобки в матричном умножении. В подсчете $s^k(y^k)^T H_k$ нужно сначала умножить $(y^k)^T H_k$, а потом вектор на вектор. Аналогично для $H_k y^k (s^k)^T$. Получается, что сложность BFGS есть $O(d^2)$ операций (не учитывая подсчет градиентов).

- Вопрос: чтобы посчитать новую матрицу нужно $O(d^2)$ операций (не учитывая подсчет градиентов). Кажется, что BFGS подсчет дороже (есть перемнножение трех матриц). Так ли это? $H_{k+1} = (I \rho_k s^k (y^k)^T) H_k (I \rho_k y^k (s^k)^T) + \rho_k s^k (s^k)^T$
- Нужно раскрыть скобки в матричном умножении. В подсчете $s^k(y^k)^T H_k$ нужно сначала умножить $(y^k)^T H_k$, а потом вектор на вектор. Аналогично для $H_k y^k (s^k)^T$. Получается, что сложность BFGS есть $O(d^2)$ операций (не учитывая подсчет градиентов).
- При инициализация матрицы H_0 достаточно брать равно единичной. Есть и более хитрые способы, но особо разницы не чувствует все работает хорошо.

• Квазиньютоновские методы не требуют подсчет гессиана и его обращение. Сложность всех арифметических операций на одной итерации $O(d^2)$, что дешевле обращения гессиана за $O(d^3)$.

- Квазиньютоновские методы не требуют подсчет гессиана и его обращение. Сложность всех арифметических операций на одной итерации $O(d^2)$, что дешевле обращения гессиана за $O(d^3)$.
- Квазиньютоновские методы имеют глобальную сверхлинейную скорость сходимости. Это медленнее, чем метод Ньютона, но не нужна «хорошая» окрестность решения.

- Квазиньютоновские методы не требуют подсчет гессиана и его обращение. Сложность всех арифметических операций на одной итерации $O(d^2)$, что дешевле обращения гессиана за $O(d^3)$.
- Квазиньютоновские методы имеют глобальную сверхлинейную скорость сходимости. Это медленнее, чем метод Ньютона, но не нужна «хорошая» окрестность решения.
- Квазиньютоновские методы используют только градиент, но в теории сходятся быстрее ускоренного градиентного метода. Вопрос: почему так, ведь метод Нестерова оптимальный?

- Квазиньютоновские методы не требуют подсчет гессиана и его обращение. Сложность всех арифметических операций на одной итерации $O(d^2)$, что дешевле обращения гессиана за $O(d^3)$.
- Квазиньютоновские методы имеют глобальную сверхлинейную скорость сходимости. Это медленнее, чем метод Ньютона, но не нужна «хорошая» окрестность решения.
- Квазиньютоновские методы используют только градиент, но в теории сходятся быстрее ускоренного градиентного метода. Вопрос: почему так, ведь метод Нестерова оптимальный? Смотри определения класса задач, для которого метод Нестерова оптимальный: не разрешены векторные произведения.

- Квазиньютоновские методы не требуют подсчет гессиана и его обращение. Сложность всех арифметических операций на одной итерации $O(d^2)$, что дешевле обращения гессиана за $O(d^3)$.
- Квазиньютоновские методы имеют глобальную сверхлинейную скорость сходимости. Это медленнее, чем метод Ньютона, но не нужна «хорошая» окрестность решения.
- Квазиньютоновские методы используют только градиент, но в теории сходятся быстрее ускоренного градиентного метода. Вопрос: почему так, ведь метод Нестерова оптимальный? Смотри определения класса задач, для которого метод Нестерова оптимальный: не разрешены векторные произведения.
- Метод Ньютона и квазиньютоновские методы на практике быстро находят точный локальный миннимум. Их спокойно можно использовать в качестве «дорешивателей». Квазиньютоновские методы и как «стартовый» метод.

Что еще можно брать вместо гессиана: $x^{k+1} = x^k - \gamma_k(B_k)^{-1} \nabla f(x^k)$?

Что еще можно брать вместо гессиана: $x^{k+1} = x^k - \gamma_k(B_k)^{-1} \nabla f(x^k)$?

• Постоянную матрицу: $B_k = B$.

Что еще можно брать вместо гессиана: $x^{k+1} = x^k - \gamma_k(B_k)^{-1} \nabla f(x^k)$?

- Постоянную матрицу: $B_k = B$.
- Аппроксимацию гессиана:

$$D_k = \operatorname{diag}\left(u_k \odot \nabla^2 f(x^k) u_k\right),$$

здесь \odot покомпонентное произведение векторов, компоненты вектора u_k берутся независимые случайные величины равные -1 и 1 с вероятностью 1/2.

Что еще можно брать вместо гессиана: $x^{k+1} = x^k - \gamma_k(B_k)^{-1} \nabla f(x^k)$?

• Постоянную матрицу: $B_k = B$.

Аппроксимацию гессиана:

$$D_k = \operatorname{diag}\left(u_k \odot \nabla^2 f(x^k) u_k\right),$$

здесь \odot покомпонентное произведение векторов, компоненты вектора u_k берутся независимые случайные величины равные -1 и 1 с вероятностью 1/2. Вопрос: что можно сказать про $\mathbb{E}D_k$?

Что еще можно брать вместо гессиана: $x^{k+1} = x^k - \gamma_k(B_k)^{-1} \nabla f(x^k)$?

- Постоянную матрицу: $B_k = B$.
- Аппроксимацию гессиана:

$$D_k = \operatorname{diag}\left(u_k \odot \nabla f(x^k)u_k\right),$$

здесь \odot покомпонентное произведение векторов, компоненты вектора u_k берутся независимые случайные величины равные -1 и 1 с вероятностью 1/2. Вопрос: что можно сказать про $\mathbb{E}D_k$? $\mathbb{E}D_k = \mathrm{diag}(\nabla^2 f(x^k))$.

Что еще можно брать вместо гессиана: $x^{k+1} = x^k - \gamma_k(B_k)^{-1} \nabla f(x^k)$?

- Постоянную матрицу: $B_k = B$.
- Аппроксимацию гессиана:

$$D_k = \operatorname{diag}\left(u_k \odot \nabla f(x^k)u_k\right),$$

здесь \odot покомпонентное произведение векторов, компоненты вектора u_k берутся независимые случайные величины равные -1 и 1 с вероятностью 1/2. Вопрос: что можно сказать про $\mathbb{E}D_k$? $\mathbb{E}D_k = \text{diag}(\nabla^2 f(x^k))$. Вопрос: хорошая ли эта аппроксимация?

Что еще можно брать вместо гессиана: $x^{k+1} = x^k - \gamma_k(B_k)^{-1} \nabla f(x^k)$?

- Постоянную матрицу: $B_k = B$.
- Аппроксимацию гессиана:

$$D_k = \operatorname{diag}\left(u_k \odot \nabla f(x^k)u_k\right),$$

здесь \odot покомпонентное произведение векторов, компоненты вектора u_k берутся независимые случайные величины равные -1 и 1 с вероятностью 1/2. Вопрос: что можно сказать про $\mathbb{E}D_k$? $\mathbb{E}D_k = \mathrm{diag}(\nabla^2 f(x^k))$. Вопрос: хорошая ли эта аппроксимация? Не особо, \mathbb{E} – это хорошо, но разброс может быть огромным.



Что еще можно брать вместо гессиана: $x^{k+1} = x^k - \gamma_k(B_k)^{-1} \nabla f(x^k)$?

- Постоянную матрицу: $B_k = B$.
- Аппроксимацию гессиана:

$$D_k = \operatorname{diag}\left(u_k \odot \nabla f(x^k)u_k\right),$$

здесь \odot покомпонентное произведение векторов, компоненты вектора u_k берутся независимые случайные величины равные -1 и 1 с вероятностью 1/2. Вопрос: что можно сказать про $\mathbb{E}D_k$? $\mathbb{E}D_k=\operatorname{diag}(\nabla^2 f(x^k))$. Вопрос: хорошая ли эта аппроксимация? Не особо, \mathbb{E} – это хорошо, но разброс может быть огромным. Поэтому делают так

$$B_{k+1} = (1 - \beta_k)B_k + \beta_k D_k,$$

где $\beta_k \in (0;1)$ (часто β_k близко к 0).



Что еще можно брать вместо гессиана: $x^{k+1} = x^k - \gamma_k(B_k)^{-1} \nabla f(x^k)$?

- Постоянную матрицу: $B_k = B$.
- Аппроксимацию гессиана:

$$D_k = \operatorname{diag}\left(u_k \odot \nabla f(x^k)u_k\right),$$

здесь \odot покомпонентное произведение векторов, компоненты вектора u_k берутся независимые случайные величины равные -1 и 1 с вероятностью 1/2. Вопрос: что можно сказать про $\mathbb{E}D_k$? $\mathbb{E}D_k=\operatorname{diag}(\nabla^2 f(x^k))$. Вопрос: хорошая ли эта аппроксимация? Не особо, \mathbb{E} – это хорошо, но разброс может быть огромным. Поэтому делают так

$$B_{k+1} = (1 - \beta_k)B_k + \beta_k D_k,$$

где $\beta_k \in (0;1)$ (часто β_k близко к 0). Такой подход помогает бороться со стохастикой. Мы аккумулируем апроксимации D_k , предполагая, что гессиан меняется слабо. С другой стороны все более старые апроксимации гессиана забываются.

Что еще можно брать вместо гессиана: $x^{k+1} = x^k - \gamma_k(B_k)^{-1} \nabla f(x^k)$?

Что еще можно брать вместо гессиана: $x^{k+1} = x^k - \gamma_k(B_k)^{-1}\nabla f(x^k)$?

• RMSProp:

$$B_{k+1}^{2} = (1-\beta)B_{k}^{2} + \beta D_{k} \quad D_{k} = \operatorname{diag}\left(f(x^{k}) \odot \nabla f(x^{k})\right)$$
• Adam:
$$\begin{pmatrix} 2 \\ \beta \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 1 \\ 1 & 1 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 1 \\ 1 & 1 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 1 \\ 1 & 1 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda \downarrow \begin{pmatrix} 2 & 0 \\ 0 & 200 \end{pmatrix} & \lambda$$

где $\beta_k = \frac{\beta - \beta^{k+1}}{1 - \beta^{k+1}}$. Лучше на начальных итерация из-за подбора β_k , который меньше «доверяет» начальным аппроксимациям.

Bonpoc: чем может помочь такого рода предобработка? подумайте о квадратичной задаче с диагональной матрицей.

Что еще можно брать вместо гессиана: $x^{k+1} = x^k - \gamma_k(B_k)^{-1} \nabla f(x^k)$?

• RMSProp:

$$B_{k+1} = (1-eta)B_k + eta D_k \quad D_k = \operatorname{diag}\left(f(x^k) \odot
abla f(x^k)
ight)$$

Adam:

$$B_{k+1} = (1 - \beta_k)B_k + \beta_k D_k \quad D_k = \operatorname{diag}\left(f(x^k) \odot \nabla f(x^k)\right),$$

где $\beta_k = \frac{\beta - \beta^{k+1}}{1 - \beta^{k+1}}$. Лучше на начальных итерация из-за подбора β_k , который меньше «доверяет» начальным аппроксимациям.

Bonpoc: чем может помочь такого рода предобработка? подумайте о квадратичной задаче с диагональной матрицей. А исчерпывающий теоретический ответ, почему условный Adam работает хорошо, не знает никто.

