

Минимизация суммы квадратов

Сергей Кривохатский

СПбГУ, Современное программирование

21 октября 2022 г.

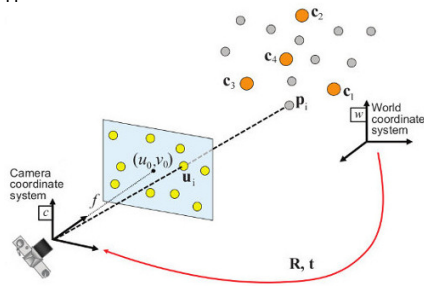
PnP

При наличии начального приближения задачу PnP можно решать путем минимизации суммы квадратов

$$\sum_i \|\pi(\mathbf{r}, \mathbf{t}, X_i) - x_i\|^2 \rightarrow \min$$

где

- ▶ \mathbf{r}, \mathbf{t} — позиция камеры
- ▶ x_i и X_i — 2D- и 3D-точки
- ▶ π — функция проекции

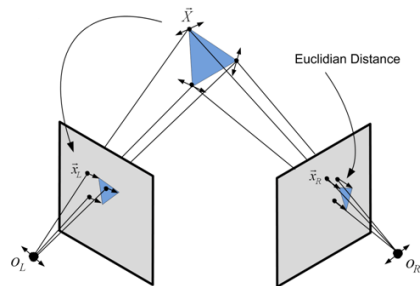


Bundle Adjustment

Качество восстановления сцены можно значительно улучшить путем оптимизации позиций камеры и координат 3D-точек

$$\sum_j \sum_{i \in V_j} \|\pi(r_j, t_j, X_i) - x_{ij}\|^2 \rightarrow \min$$

- ▶ r_j, t_j — позиция камеры в кадре j
- ▶ X_i — 3D-позиция i -й точки
- ▶ V_j — номера видимых в кадре j точек
- ▶ x_{ij} — 2D-позиция i -й точки в кадре j
- ▶ π — функция проекции



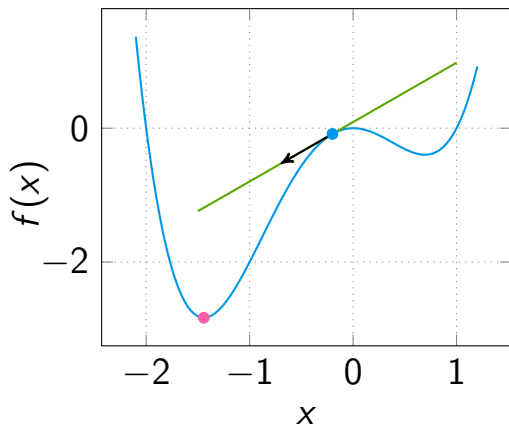
Оптимизация на примере одномерной функции

Задача оптимизации

- ▶ дана $f: \mathbb{R} \rightarrow \mathbb{R}$
- ▶ хотим найти
 $x_{\min} = \arg \min_x f(x)$
- ▶ знаем начальное
приближение x_0

Градиентный спуск

- ▶ повторяем
 $x_{k+1} = x_k - \alpha_{k+1} f'(x_k)$
- ▶ пока $|f'(x_k)| \geq \varepsilon$ (например)

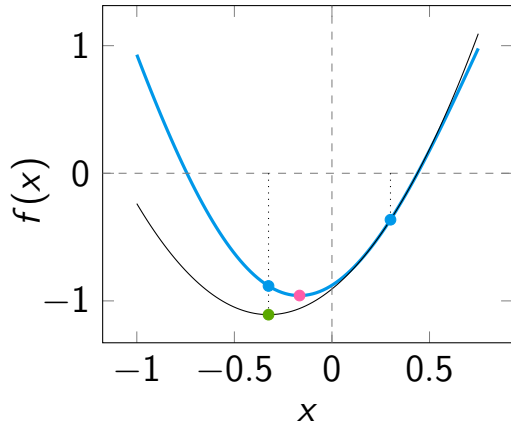


Метод Ньютона для одномерной функции

Аппроксимируем $f(x)$
параболами и ходим
в минимумы

$$f(x_k + \delta) \approx f(x_k) + f'(x_k)\delta + \frac{1}{2}f''(x_k)\delta^2$$

При определенных условиях
скорость сходимости
квадратичная



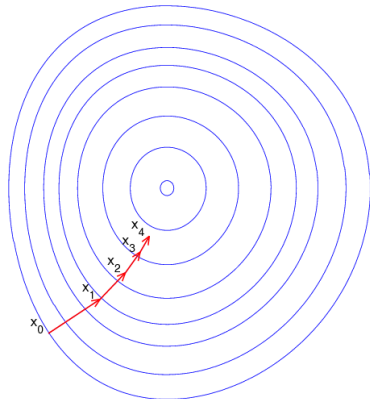
Оптимизация многомерной функции

Задача оптимизации

- ▶ дана $f: \mathbb{R}^n \rightarrow \mathbb{R}$
- ▶ хотим найти
$$x_{\min} = \arg \min_x f(x)$$
- ▶ знаем начальное приближение x_0

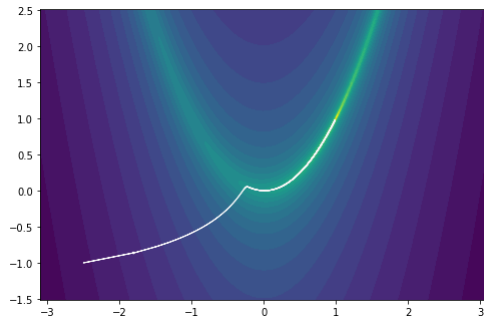
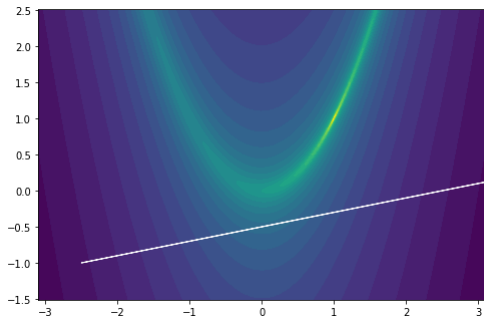
Градиентный спуск

- ▶ повторяем
$$x_{k+1} = x_k - \alpha_{k+1} \nabla f(x_k)$$
- ▶ пока $\|\nabla f(x_k)\| \geq \varepsilon$ (например)



Градиентный спуск

- ▶ обычно требует много коротких шагов, чтобы сойтись
- ▶ стоимость вычисления одного шага относительно небольшая



Метод Ньютона для многомерной функции I

- ▶ известна позиция x_k
- ▶ приближаем f в точке x_k квадратичной функцией

$$f(x_k + \delta) \approx f(x_k) + \delta^T \nabla f(x_k) + \frac{1}{2} \delta^T \mathbf{H}_f(x_k) \delta =: f_k(\delta)$$

- ▶ вычисляем $\delta_k = \arg \min_{\delta} f_k(\delta)$
- ▶ следующая позиция: $x_{k+1} = x_k + \delta_k$
- ▶ повторяем, пока не выполнится условие остановки

Метод Ньютона для многомерной функции II

$$f_k(\delta) := f(x_k) + \delta^\top \nabla f(x_k) + \frac{1}{2} \delta^\top \mathbf{H}_f(x_k) \delta$$

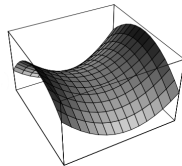
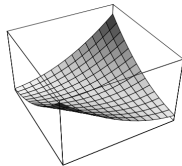
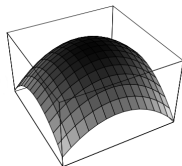
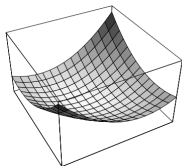
► $\nabla f(x_k) = \left(\frac{\partial f(x_k)}{\partial x_1} \quad \dots \quad \frac{\partial f(x_k)}{\partial x_n} \right)^\top =: \mathbf{g}_k$ — градиент

► $\mathbf{H}_f(x_k) = \begin{pmatrix} \frac{\partial^2 f(x_k)}{\partial x_1 \partial x_1} & \dots & \frac{\partial^2 f(x_k)}{\partial x_1 \partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f(x_k)}{\partial x_n \partial x_1} & \dots & \frac{\partial^2 f(x_k)}{\partial x_n \partial x_n} \end{pmatrix} =: \mathbf{H}_k$ — матрица Гессе

Метод Ньютона для многомерной функции III

$$f_k(\delta) := f(x_k) + \delta^T g_k + \frac{1}{2} \delta^T H_k \delta$$

- ▶ f_k имеет единственный минимум, если H_k — положительно определенная
- ▶ чтобы его найти, нужно решить СЛАУ: $H_k \delta = -g_k$
- ▶ если H_k не положительно определенная, нам не повезло



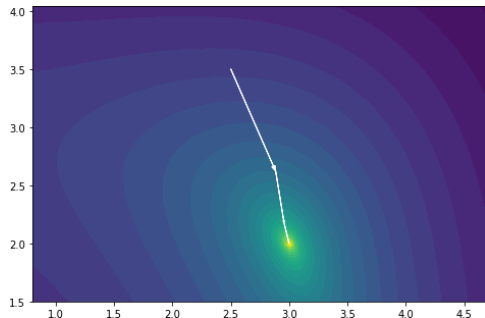
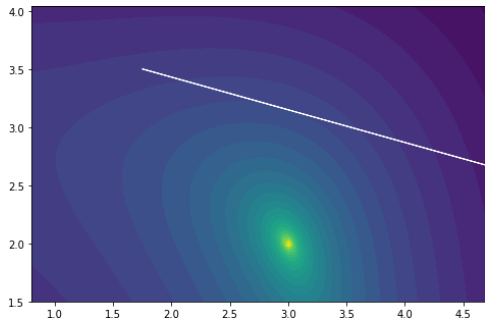
Метод Ньютона для многомерной функции IV

Получаем следующий алгоритм

- ▶ известна позиция x_k
- ▶ приближаем f в точке x_k квадратичной функцией $f_k(\delta)$
- ▶ вычисляем шаг, на котором f_k достигает минимума
$$\delta_k = -\mathbf{H}_k^{-1} \mathbf{g}_k$$
- ▶ следующая позиция: $x_{k+1} = x_k + \delta_k$
- ▶ повторяем, пока не выполнится условие остановки (какое-нибудь)

Метод Ньютона для многомерной функции V

- ▶ очень чувствителен к начальному приближению
- ▶ если все хорошо, очень быстро сходится
- ▶ стоимость вычисления одного шага относительно велика



Метод Гаусса — Ньютона I

- ▶ вычислять матрицу Гессе напрямую может быть сложно
- ▶ будем делать это приближенно
- ▶ мы знаем, что f имеет особый вид

$$f(x) = \sum_{i=1}^m r_i^2(x) = r(x)^T r(x)$$

где $r(x) = (r_1(x) \ r_2(x) \ \cdots \ r_m(x))^T$

Метод Гаусса — Ньютона II

Распишем элемент матрицы Гессе

$$\frac{\partial^2 r^\top r}{\partial x_i \partial x_j} = \frac{\partial}{\partial x_i} \left(2r^\top \frac{\partial r}{\partial x_j} \right) = 2 \frac{\partial r^\top}{\partial x_i} \frac{\partial r}{\partial x_j} + 2r^\top \frac{\partial^2 r}{\partial x_i \partial x_j}$$

Если элементы $r(x)$ достаточно малы либо близки к линейным функциям (т. е. их вторые производные малы)

$$\frac{\partial^2 f}{\partial x_i \partial x_j} \approx 2 \frac{\partial r^\top}{\partial x_i} \frac{\partial r}{\partial x_j}$$

Метод Гаусса — Ньютона III

Матрица Якоби функции $r(x)$

$$\begin{aligned} J_r(x) &= \begin{pmatrix} \frac{\partial r(x)}{\partial x_1} & \dots & \frac{\partial r(x)}{\partial x_n} \end{pmatrix} = \\ &= \begin{pmatrix} \frac{\partial r_1(x)}{\partial x_1} & \dots & \frac{\partial r_1(x)}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial r_m(x)}{\partial x_1} & \dots & \frac{\partial r_m(x)}{\partial x_n} \end{pmatrix} \end{aligned}$$

Несложно заметить, что

$$H_f(x) \approx 2J_r(x)^T J_r(x)$$

И что

$$\nabla f(x) = 2J_r(x)^T r(x)$$

так как

$$\frac{\partial f}{\partial x_i} = 2 \frac{\partial r}{\partial x_i}^T r$$

Метод Гаусса — Ньютона IV

- ▶ обозначим $r_k := r(x_k)$ и $J_k := J_r(x_k)$
- ▶ f в точке x_k приближается как

$$\begin{aligned} f_k(\delta) &\approx f(x_k) + \delta^\top g_k + \frac{1}{2} \delta^\top H_k \delta \approx \\ &\approx r_k^\top r_k + \delta^\top 2J_k^\top r_k + \delta^\top J_k^\top J_k \delta \end{aligned}$$

- ▶ шаг δ_k получается из решения линейной системы

$$J_k^\top J_k \delta_k = -J_k^\top r_k$$

Метод Гаусса — Ньютона V

Получаем следующий алгоритм

- ▶ известна позиция x_k
- ▶ приближаем f в точке x_k квадратичной функцией:
$$f_k(\delta) \approx r_k^T r_k + \delta^T 2J_k^T r_k + \delta^T J_k^T J_k \delta$$
- ▶ вычисляем шаг, минимизируя приближение:
$$\delta_k = -(J_k^T J_k)^{-1} J_k^T r_k$$
- ▶ следующая позиция: $x_{k+1} = x_k + \delta_k$
- ▶ повторяем, пока не выполнится условие остановки (какое-нибудь)

Алгоритм Левенберга — Марквардта I

Матрица $\mathbf{J}_k^T \mathbf{J}_k$ может быть вырожденной. Левенберг в 1944 году предложил использовать регуляризацию $\mathbf{J}_k^T \mathbf{J}_k + \lambda \mathbf{1}$

$$f_k(\delta) = f(x_k) + \delta^T g_k + \delta^T \mathbf{J}_k^T \mathbf{J}_k \delta + \delta^T \lambda \mathbf{1} \delta$$

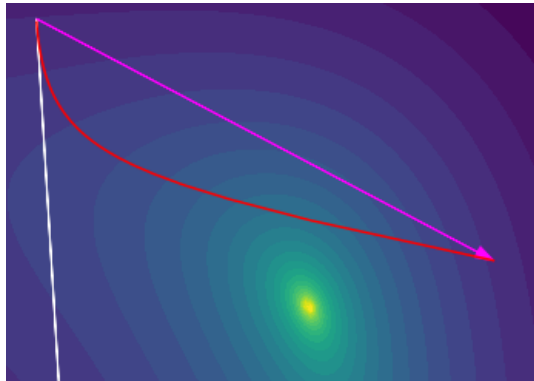
Прибавляем параболоид с минимумом в x_k



Алгоритм Левенберга — Марквардта II

$$\delta_k = -(\mathbf{J}_k^T \mathbf{J}_k + \lambda \mathbf{1})^{-1} \mathbf{g}_k$$

- ▶ меньше λ — метод Гаусса — Ньютона
- ▶ больше λ — градиентный спуск с шагом $1/\lambda$



Алгоритм Левенберга — Марквардта III

- ▶ Марквардт в 1963 году предложил менять λ адаптивно от шага к шагу по следующему принципу

$$\lambda_{k+1} = \begin{cases} \frac{1}{\nu} \lambda_k & \text{если } f(x_k + \delta_k) < f(x_k) \\ \nu \lambda_k, & \text{иначе} \end{cases}$$

где $\nu > 1$ — константа (часто берут 10)

- ▶ существует формулировка алгоритма подбора λ в терминах доверительного региона, которая обладает некоторыми преимуществами, см. *J. More, "The Levenberg-Marquardt Algorithm: Implementation and Theory", 1977*

Алгоритм Левенберга — Марквардта IV

- ▶ Флетчер в 1971 году предложил вместо **1** брать $\text{diag } \mathbf{J}_k^T \mathbf{J}_k$ — это улучшает сходимость в направлениях с малым абсолютным значением градиента
- ▶ может возникнуть проблема: если параметр x_i не влияет на функцию, матрица $\mathbf{J}_k^T \mathbf{J}_k + \lambda_k \text{diag } \mathbf{J}_k^T \mathbf{J}_k$ становится вырожденной

Алгоритм Левенберга — Марквардта V

Шаг алгоритма Левенберга — Марквардта

$$\delta_k = -(\mathbf{J}_k^T \mathbf{J}_k + \lambda_k \text{diag } \mathbf{J}_k^T \mathbf{J}_k)^{-1} \mathbf{J}_k^T r_k$$

Вычислять обратную матрицу не стоит, лучше решить СЛАУ

$$(\mathbf{J}_k^T \mathbf{J}_k + \lambda_k \text{diag } \mathbf{J}_k^T \mathbf{J}_k) \delta_k = -\mathbf{J}_k^T r_k$$

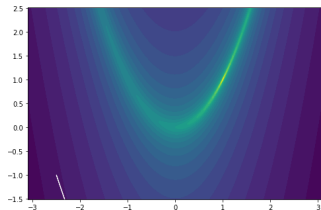
Матрица \mathbf{J}_k часто имеет специфическую структуру, что может позволить подобрать наиболее эффективный метод решения, оптимизировать потребление памяти и т. д.

Алгоритм Левенберга — Марквардта VI

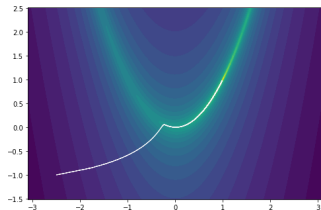
Получаем следующую схему алгоритма

- ▶ известна позиция x_k
- ▶ вычисляем шаг, решая систему уравнений
$$(\mathbf{J}_k^T \mathbf{J}_k + \lambda_k \text{diag } \mathbf{J}_k^T \mathbf{J}_k) \delta_k = -\mathbf{J}_k^T r_k$$
- ▶ если $f(x_{k+1} + \delta_k) \geq f(x_k)$, увеличиваем λ , следующая позиция: $x_{k+1} = x_k$
- ▶ иначе уменьшаем значение λ , следующая позиция: $x_{k+1} = x_k + \delta_k$
- ▶ повторяем, пока не выполнится условие остановки (какое-нибудь)

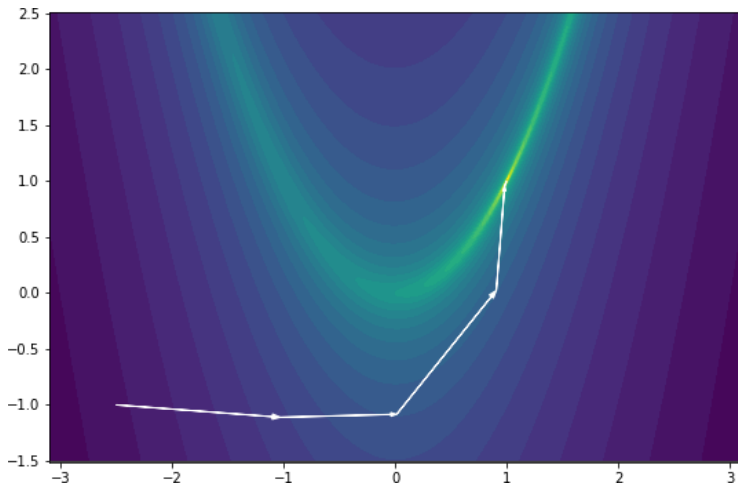
Алгоритм Левенберга — Марквардта VII



Метод Гаусса — Ньютона



Градиентный спуск



Метод Левенберга — Марквардта

Робастное оценивание на примере PnP I

- ▶ решение задачи PnP путем оптимизации

$$\arg \min_{r,t} \sum_i \|\pi(r, t, X_i) - x_i\|^2$$

- ▶ для удобства изменим обозначения

$$\sum_i \|\pi(r, t, X_i) - x_i\|^2 = \sum_i r_i^2(p) = r(p)^T r(p)$$

- ▶ предположим, что при известной позиции камеры ошибки репроекции распределены нормально и независимо

$$r_i(p) \sim \mathcal{N}(0, \sigma)$$

Робастное оценивание на примере PnP II

- ▶ тогда позицию камеры можно найти с помощью метода максимального правдоподобия

$$\arg \max_p \prod_i \frac{1}{\sigma \sqrt{2\pi}} \exp \left(-\frac{r_i(p)^2}{2\sigma^2} \right)$$

- ▶ если прологарифмировать, умножить на -1 и избавиться от констант, получим сумму квадратов

$$\arg \min_p \sum_i r_i^2(p)$$

Робастное оценивание на примере Р n Р III

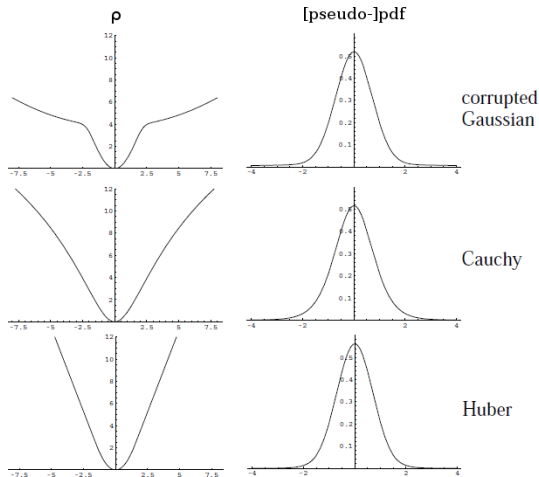
- ▶ выбросы среди 2D-точек могут заметно исказить результат
- ▶ для решения этой проблемы можно заменить квадратичную функцию на другую, которая даст более робастную оценку

$$\arg \min_p \sum_i \rho(r_i(p))$$

- ▶ такие оценки называются *оценками типа максимального правдоподобия* или *M-оценками*

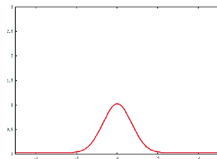
M-оценки I

- ρ подбирается исходя из предположений о распределении ошибок измерений (2D-точек) или из эвристических соображений
- ρ может соответствовать «честной» плотности распределения измерений или «нечестной» (интеграл не равен единице, например)



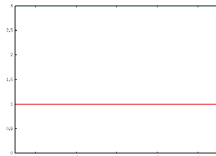
M-оценки II

Normal distribution
(inliers)



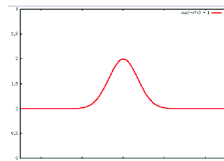
+

Uniform distribution
(outliers)

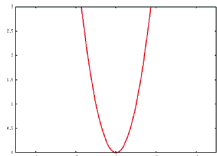


=

Mixture

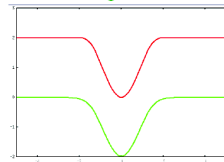


$-\log(\cdot)$



Least-squares

$-\log(\cdot)$



Tukey estimator

М-оценки III

- ▶ какие-то М-оценки (например, Тьюки) являются невыпуклыми, из-за чего могут возникнуть проблемы со сходимостью к глобальному минимуму, если начальное приближение плохое
- ▶ другие М-оценки (например, Хьюбера) являются выпуклыми и не имеют такого негативного влияния на сходимость
- ▶ на практике иногда полезно прогнать оптимизацию несколько раз: например, сначала с М-оценкой Хьюбера, потом с Тьюки
- ▶ *A General and Adaptive Robust Loss Function (CVPR 2019)*

M-оценки. Оптимизация I

- ▶ ρ — M-оценка
- ▶ $\psi := \rho'$ — функция влияния
- ▶ $w(e) := \frac{\psi(e)}{e}$ — весовая функция

$$f(p) := r(p)^T r(p)$$

$$f_\rho(p) := \sum_{k=0}^m \rho(r_k(p))$$

M-оценки. Оптимизация II

Продифференцируем $\rho(r_k(p))$ по p_i

$$\begin{aligned}\frac{\partial \rho(r_k(p))}{\partial p_i} &= \psi(r_k(p)) \frac{\partial r_k(p)}{\partial p_i} = \\ &= w(r_k(p)) r_k(p) \frac{\partial r_k(p)}{\partial p_i}\end{aligned}$$

$$\psi := \rho'$$

$$w(e) := \frac{\psi(e)}{e}$$

Значит градиент можно выразить таким образом

$$\nabla f_\rho(p) = J_r^\top(p) W_p r(p)$$

(W_p — диагональная матрица с весами в точке p)

М-оценки. Оптимизация III

Продифференцируем еще раз, теперь по p_j

$$\frac{\partial}{\partial p_j} \left(\psi(r_k(p)) \frac{\partial r_k(p)}{\partial p_i} \right) = \frac{\partial r_k(p)}{\partial p_j} \psi'(r_k(p)) \frac{\partial r_k(p)}{\partial p_i} + \psi(r_k(p)) \frac{\partial^2 r_k(p)}{\partial p_i \partial p_j}$$

Значит матрицу Гессе можно приблизить таким образом

$$\mathbf{H}_{f_p}(p) \approx \mathbf{J}_r^T(p) \operatorname{diag}(\psi'(r_1(p)), \psi'(r_2(p)), \dots) \mathbf{J}_r(p)$$

Но на практике обычно делают даже так

$$\mathbf{H}_{f_p}(p) \approx \mathbf{J}_r^T(p) \mathbf{W}_p \mathbf{J}_r(p)$$

М-оценки. Оптимизация IV

- ▶ минимизацию $\sum \rho(r_k(p))$ сводят к шагам с минимизацией взвешенных сумм квадратов $r(p)^T \mathbf{W}_p r(p)$
- ▶ на каждом шаге оптимизации производят перевзвешивание, т. е. пересчитывают весовую матрицу \mathbf{W}
- ▶ на практике параметры М-оценки (т. е. точная форма функции ρ) обычно переоцениваются на каждом шаге на основе текущего распределения r_k

Спасибо за внимание!