

Э. РАЙК

ДИФФЕРЕНЦИРУЕМОСТЬ ПО ПАРАМЕТРУ ФУНКЦИИ ВЕРОЯТНОСТИ И СТОХАСТИЧЕСКИЙ ПСЕВДОГРАДИЕНТНЫЙ МЕТОД ДЛЯ ЕЕ ОПТИМИЗАЦИИ

Пусть на n -мерном пространстве X задана функция вероятности

$$v(x) = P\{f(x, \xi(\omega)) \geq 0\}, \quad (1)$$

где $\xi(\omega)$ — случайный вектор со значениями в m -мерном пространстве Y , $f(x, y)$ — действительная функция, определенная на произведении пространств $X \times Y$.

Допустим ниже, что случайный вектор $\xi(\omega)$ индуцирует в пространстве Y вероятностное распределение, которое имеет лишь абсолютно непрерывное составляющее и, следовательно, полностью определяется своей функцией плотности вероятности $p(y)$. Это позволяет функцию вероятности выписать через m -кратный интеграл

$$v(x) = \int_{f(x,y) \geq 0} \dots \int p(y) dy. \quad (2)$$

При решении задач нелинейного стохастического программирования, содержащих функции вероятности в качестве функций ограничения или функции цели, встает вопрос о виде градиента и условий дифференцируемости функции вероятности.

Ответ на этот вопрос дает

Теорема 1. Пусть 1) частные производные $f'_x(x, y)$, $f'_y(x, y)$ и плотность вероятности $p(y)$ непрерывны;

2) множество $S(x, 0) = \{y : f(x, y) = 0\}$ равномерно ограничено в некоторой окрестности точки x ;

3) функция $f'_y(x, y)$ удовлетворяет условию Липшица по y и функция $\|f'_y(x, y)\|$ положительна на множестве $S(x, 0)$;

4) поверхностный интеграл $\int_{S(x,0)} \dots \int dS$ существует.

Тогда функция вероятности $v(x)$ дифференцируема по Фреше в точке x и градиент $v'(x)$ равняется $m - 1$ -кратному интегралу

$$v'(x) = \int_{S(x,0)} \dots \int f'_x(x, y) p(y) \|f'_y(x, y)\|^{-1} dS. \quad (3)$$

Доказательство. В силу предположения 2 и непрерывности функции $f(x, y)$ поверхность $S(x + \Delta x, 0) = \{y : f(x + \Delta x, y) = 0\}$ приближается к поверхности $S(x, 0)$ равномерно в том смысле, что максимальное расстояние $\varrho(\Delta x)$ точек поверхности $S(x + \Delta x, 0)$ до поверхности $S(x, 0)$ стремится к нулю, т. е. $\varrho(\Delta x) \rightarrow 0$ при $\|\Delta x\| \rightarrow 0$. Действительно, так как множества $S(x, 0)$ и $S(x + \Delta x, 0)$ компактны, расстояние $\varrho(\Delta x)$ существует и равняется

$$\varrho(\Delta x) = \max_{u \in S(x+\Delta x, 0)} \min_{w \in S(x, 0)} \|u - w\|.$$

Предположим обратное: $\varrho(\Delta x) \not\rightarrow 0$ при $\|\Delta x\| \rightarrow 0$. Тогда найдется такая сходящаяся последовательность точек $\{x_n, y_n\}$, реализующих расстояние

$$\varrho(x_n - x) = \min_{w \in S(x, 0)} \|y_n - w\|, \quad y_n \in S(x_n, 0),$$

что $f(x_n, y_n) = 0$ и по непрерывности $f(x, y) = 0$, но это противоречит предположению $\varrho(\Delta x) \not\rightarrow 0$.

При непрерывных частных производных $f'_x(x, y)$ и $f'_y(x, y)$ имеем разложение

$$f(x + \Delta x, y + \Delta y) = f(x, y) + (f'_x(\tilde{x}, \tilde{y}), \Delta x) + (f'_y(\tilde{x}, \tilde{y}), \Delta y), \quad (4)$$

где $\tilde{x} = x + \Theta(y)\Delta x$, $\tilde{y} = y + \Theta(y)\Delta y$, $0 \leq \Theta(y) \leq 1$.

Здесь через (x, y) обозначено скалярное произведение векторов x и y . Функцию $\Theta(y)$ можно считать измеримой ([1], с. 409). Приравняем в формуле (4) $f(x, y) = 0$ и $f(x + \Delta x, y + \Delta y) = 0$ и выразим из формулы (4) приращение Δy , взятое параллельно вектору нормали поверхности $S(x, 0)$ в точке y , через

$$\Delta y = - \frac{f'_y(x, y) (f'_x(\tilde{x}, \tilde{y}), \Delta x)}{(f'_y(\tilde{x}, \tilde{y}), f'_y(x, y))}.$$

В силу условия 3 теоремы существует Δx такая, начиная с которой $(f'_y(\tilde{x}, \tilde{y}), f'_y(x, y)) \geq \gamma > 0$ для всех y поверхности $S(x, 0)$. Проекция приращения Δy на внешнюю нормаль $-f'_y(x, y) \|f'_y(x, y)\|^{-1}$ множества $\{y : f(x, y) \geq 0\}$ равняется

$$l(y) = \frac{(f'_x(\tilde{x}, \tilde{y}), \Delta x) \|f'_y(x, y)\|}{(f'_y(\tilde{x}, \tilde{y}), f'_y(x, y))}$$

Функция $l(y)$ имеет вполне определенный геометрический смысл. Она равняется по абсолютной величине отрезку между поверхностями $S(x, 0)$ и $S(x + \Delta x, 0)$, взятому параллельно вектору нормали поверхности $S(x, 0)$. Функция $l(y)$ является положительной, если отрезок расположен вне множества $\{y : f(x, y) \geq 0\}$, и отрицательной, если отрезок расположен внутри множества.

Приращение функции вероятности $v(x)$ выражается в виде разности m -кратных интегралов от функции $p(y)$

$$v(x + \Delta x) - v(x) = \int_{V_1} \dots \int p(y) dy - \int_{V_2} \dots \int p(y) dy,$$

где $V_1 = \{y : f(x + \Delta x, y) \geq 0, f(x, y) < 0\}$

и

$$V_2 = \{y : f(x + \Delta x, y) < 0, f(x, y) \geq 0\}.$$

Эти m -кратные интегралы можно выразить через $m - 1$ -кратный поверхностный интеграл

$$v(x + \Delta x) - v(x) = \int_{S(x, 0)} \dots \int l(y) p(h(y)) dS + o(\|\Delta x\|),$$

где $h(y) = y + f'_y(x, y) \|f'_y(x, y)\|^{-1} l(y) \Theta_1(y)$, $0 \leq \Theta_1(y) \leq 1$.

Это выражение можно получить непосредственно, учитывая то, что объем элементарного кривого цилиндра ΔV , боковая поверхность которого образуется нормальными поверхностями $S(x, 0)$, с точностью до бесконечно малых высшего порядка относительно $\|\Delta x\|$, равняется $\Delta V = \Delta S_x l(\bar{y})$, $\bar{y} \in S_x$, где ΔS_x — площадь основания элементарного цилиндра.

По определению функция $l(y)$ положительна для области V_1 и отрицательна для области V_2 .

Итак,

$$v(x+\Delta x) - v(x) = \int_{S(x,0)} \dots \int (f'_x(x,y), \Delta x) p(y) \|f'_y(x,y)\|^{-1} dS + \omega(x, \Delta x),$$

причем остаточный член равняется

$$\begin{aligned} \omega(x, \Delta x) &= \int_{S(x,0)} \dots \int \frac{(f'_x(\tilde{x}, \tilde{y}), \Delta x) p(h(y)) \|f'_y(x,y)\|}{(f'_y(\tilde{x}, \tilde{y}), f'_y(x,y))} dS - \\ &- \int_{S(x,0)} \dots \int \frac{(f'_x(x,y), \Delta x) p(y)}{\|f'_y(x,y)\|} dS + o(\|\Delta x\|), \end{aligned}$$

$$\begin{aligned} \tilde{x} &= x + \Theta(y) \Delta x, & h(y) &= y + \frac{f'_y(x,y) (f'_x(\tilde{x}, \tilde{y}), \Delta x)}{(f'_y(\tilde{x}, \tilde{y}), f'_y(x,y))} \Theta_1(y), \\ \tilde{y} &= y + \Theta(y) \Delta y, & & \\ 0 \leq \Theta(y) \leq 1, & & 0 \leq \Theta_1(y) \leq 1. \end{aligned}$$

Покажем, что выбирая $\|\Delta x\|$ малой, отношение $\frac{|\omega(x, \Delta x)|}{\|\Delta x\|}$ можно сделать сколь угодно малым. Допустим, что $\|\Delta x\|$ настолько мала, что $(f'_y(\tilde{x}, \tilde{y}), f'_y(x,y)) \geq \gamma > 0$ на множестве $Q = \{y : \|z - y\| \leq \rho(\Delta x), z \in S(x,0)\}$. Обозначим предельные значения непрерывных функций на компакте через

$$\max_{y \in Q} (\|f'_x(x,y)\|, \|f'_y(x,y)\|, p(y)) = K_1, \quad \min_{y \in Q} \|f'_y(x,y)\| = K_2.$$

Непрерывные функции $f'_x(x,y)$, $\|f'_y(x,y)\|$, $p(y)$ на компакте Q равномерно непрерывны, следовательно, для любого $\varepsilon > 0$ существует $\delta > 0$ такое, что если выбрать $\|\Delta x\|$ настолько малой, что $\rho(\Delta x) \leq \delta$, то

$$\max_{\|\tilde{y}-y\| \leq \|\Delta y\|} (\|f'_x(\tilde{x}, \tilde{y}) - f'_x(x,y)\|, \|f'_y(\tilde{x}, \tilde{y})\| - \|f'_y(x,y)\|, |p(h(y)) - p(y)|) \leq \varepsilon.$$

Тогда остаточный член оценивается таким образом

$$\begin{aligned} \frac{|\omega(x, \Delta x)|}{\|\Delta x\|} &\leq \int_{S(x,0)} \dots \int \left\| \frac{f'_x(\tilde{x}, \tilde{y}) p(h(y))}{\|f'_y(\tilde{x}, \tilde{y})\|} - \frac{f'_x(x,y) p(y)}{\|f'_y(x,y)\|} \right\| dS \leq \\ &\leq \int_{S(x,0)} \dots \int dS \frac{3K_1^2}{K_2^2} \varepsilon. \end{aligned}$$

По предположению теоремы поверхностный интеграл $\int_{S(x,0)} \dots \int dS$ существует. Теорема доказана.

Наиболее жестким из условий теоремы 1 является требование ограниченности множества $S(x,0) = \{y : f(x,y) = 0\}$. Если отказаться от этого требования, то верна

Теорема 2. Пусть 1) частная производная $f'_x(x,y)$ и плотность вероятности $p(y)$ непрерывны;

2) функция $f'_y(x,y)$ удовлетворяет условию Липшица по y и функция $\|f'_y(x,y)\|$ положительна на множестве $S(x,0)$;

3) существует поверхностный интеграл

$$\int_{S(x,0)} \dots \int f'_x(x,y) p(y) \|f'_y(x,y)\|^{-1} dS;$$

4) для любого ограниченного множества Q существует интеграл $\int_{S(x,0) \cap Q} \dots \int dS$;

5) для любого $\varepsilon > 0$ существует ограниченное множество $M \subset Y$ и $\delta > 0$ такие, что если $\|\Delta x\| \leq \delta$, то $|v_{Y \setminus M}(x + \Delta x) - v_{Y \setminus M}(x)| \leq \varepsilon \|\Delta x\|$, где $v_{Y \setminus M}(x) = \iint_{Y \setminus M \cap \{y: f(x,y)=0\}} \dots \int p(y) dy$.

Тогда функция вероятности $v(x)$ дифференцируема по Фреше в точке x и градиент $v'(x)$ равняется $t - 1$ -кратному интегралу

$$v'(x) = \int_{S(x,0)} \dots \int f'_x(x, y) p(y) \|f'_y(x, y)\|^{-1} dS.$$

Доказательство. Зафиксируем $\varepsilon > 0$. Тогда существует ограниченное множество M и $\delta_1 > 0$ такие, что если $\|\Delta x\| \leq \delta_1$, то

$$|v_{Y \setminus M}(x + \Delta x) - v_{Y \setminus M}(x)| = |P\{f(x + \Delta x, \xi(\omega)) \geq 0, \xi(\omega) \in Y \setminus M\} - P\{f(x, \xi(\omega)) \geq 0, \xi(\omega) \in Y \setminus M\}| \leq \varepsilon \|\Delta x\|$$

$$\text{и} \quad \left| \int_{S(x,0) \cap Y \setminus M} \dots \int f'_x(x, y) p(y) \|f'_y(x, y)\|^{-1} dS \right| \leq \varepsilon.$$

Имеем

$$\left| v_{Y \setminus M}(x + \Delta x) - v_{Y \setminus M}(x) - \int_{S(x,0) \cap Y \setminus M} \dots \int (f'_x(x, y), \Delta x) p(y) \|f'_y(x, y)\|^{-1} dS \right| \leq 2\varepsilon \|\Delta x\|.$$

Согласно теореме 1, на множестве M можем выбрать $\delta_2 > 0$ такое, что при $\|\Delta x\| \leq \delta_2$

$$|v_M(x + \Delta x) - v_M(x) - \int_{S(x,0) \cap M} \dots \int (f'_x(x, y), \Delta x) p(y) \|f'_y(x, y)\|^{-1} dS| \leq \varepsilon \|\Delta x\|.$$

Итак, для любого $\varepsilon > 0$ существует $\delta = \min(\delta_1, \delta_2)$ такое, что

$$|v(x + \Delta x) - v(x) - \int_{S(x,0)} \dots \int (f'_x(x, y), \Delta x) p(y) \|f'_y(x, y)\|^{-1} dS| \leq 3\varepsilon \|\Delta x\|.$$

Рассмотрим задачу минимизации функции вероятности

$$\min_{x \in X} v(x).$$

Использование градиентного метода для решения этой задачи затруднительно, так как приведенная формула (3) для вычисления градиента функции вероятности влечет за собой громоздкие вычисления. Помимо вычисления сложного интеграла по $t - 1$ -мерной поверхности, нам потребуется еще точно знать явный вид функции плотности вероятности $p(y)$. Оказывается, что эти трудности можно обойти.

Предлагаемый ниже стохастический псевдоградиентный метод для минимизации по параметру функции вероятности $v(x) = P[f(x, \xi(\omega)) \geq 0]$ не требует ни знания плотности вероятности $p(y)$, ни проведения громоздких вычислений. В этом смысле он напоминает известный градиентный алгоритм стохастической аппроксимации Роббинса—Монро [2].

Алгоритмы стохастической аппроксимации для решения задач, содержащих функции математического ожидания, нашли широкое использование начиная с работ [2, 3]; условия их сходимости хорошо изучены. В работе [4] исследована сходимость целого класса алгоритмов стохастической оптимизации, так наз. стохастических псевдоградиентных алгоритмов.

Приведем общий вид стохастических алгоритмов

$$x_h = x_{h-1} - \gamma_h s_h, \quad (5)$$

где $x_h \in X$, γ_h — детерминированные скалярные множители, зависящие

только от k , s_k — реализации случайного вектора в X . Следуя работе [4], назовем стохастический алгоритм для функции $v(x)$ псевдоградиентным, если

$$(v'(x_k), Es_k) \geq 0, \tag{6}$$

где Es_k — математическое ожидание случайного вектора s_k .

Для функции вероятности $v(x)$ определим s_k следующим образом:

$$s_k = f'_x(x_k, y_k), \tag{7}$$

где y_k — первая очередная реализация случайного вектора, для которой

$$|f(x_k, y_k)| \leq \varepsilon_k, \quad \varepsilon_k > 0. \tag{8}$$

Смысл и выбор значений ε_k уточняется ниже.

Теорема 3. Пусть выполнены условия теоремы 1. Тогда существует последовательность чисел $\delta_k > 0$ такая, что если $\varepsilon_k \leq \delta_k$, то алгоритм (5), (7), (8) является псевдоградиентным.

Доказательство. Пусть $v'(x_k) \neq 0$. В противном случае неравенство (6) выполняется автоматически. Используя определение s_k , выразим математическое ожидание Es_k через m -кратный интеграл

$$Es_k = \iint \dots \int_{|f(x_k, y)| \leq \varepsilon_k} f'_x(x_k, y) p_k(y) dy. \tag{9}$$

Здесь $p_k(y)$ определяется через $p(y)$ с учетом того, что используется лишь та реализация y_k , для которой выполняется неравенство (8):

$$p_k(y) = \frac{p(y)}{P(x_k, \varepsilon_k)} \quad \text{и} \quad p_k(y) \equiv 0, \quad \text{если} \quad P(x_k, \varepsilon_k) = 0,$$

где $P(x_k, \varepsilon_k) = \iint \dots \int_{|f(x_k, y)| \leq \varepsilon_k} p(y) dy$.

Используя теорему 108 [5], перепишем m -кратный интеграл (9) в виде

$$\iint \dots \int_{|f(x_k, y)| \leq \varepsilon_k} f'_x(x_k, y) p_k(y) dy = \int_{-\varepsilon_k}^{\varepsilon_k} d\omega \int_{S(x_k, \omega)} f'_x(x_k, y) p_k(y) \|f'_y(x_k, y)\|^{-1} dS,$$

где $S(x_k, \omega) = \{y : f(x_k, y) = \omega\}$.

По теореме о среднем значении имеем

$$Es_k = 2\varepsilon_k \int_{S(x_k, \bar{\omega})} f'_x(x_k, y) p_k(y) \|f'_y(x_k, y)\|^{-1} dS, \quad -\varepsilon_k \leq \bar{\omega} \leq \varepsilon_k$$

и аналогично получаем

$$P(x_k, \varepsilon_k) = 2\varepsilon_k \int_{S(x_k, \tilde{\omega})} p(y) \|f'_y(x_k, y)\|^{-1} dS, \quad -\varepsilon_k \leq \tilde{\omega} \leq \varepsilon_k.$$

Так как непрерывные функции $f'_x(x_k, y)$, $\|f'_y(x_k, y)\|$, $p_k(y)$ на компакте равномерно непрерывны и поверхности $S(x_k, \bar{\omega})$ и $S(x_k, \tilde{\omega})$ приближаются к поверхности $S(x_k, 0)$ равномерно, то

$$Es_k = 2\varepsilon_k \int_{S(x_k, 0)} f'_x(x_k, y) p_k(y) \|f'_y(x_k, y)\|^{-1} dS + o_1(\varepsilon_k),$$

где $o_1(\varepsilon_k)$ — бесконечно малая величина высшего порядка относительно ε_k и

$$P(x_k, \varepsilon_k) = 2\varepsilon_k P_k + o_2(\varepsilon_k),$$

$$P_h = \int \dots \int_{S(x_k, 0)} p(y) \|f'_y(x_h, y)\|^{-1} dS.$$

Заметим, что $P_h > 0$, так как $\|f'_y(x_h, y)\| > 0$ и по предположению, сделанному вначале, $v'(x_h) \neq 0$.

Перепишем Es_h в форме

$$Es_h = \frac{1}{P_h} \int \dots \int_{S(x_k, 0)} f'_x(x_h, y) p(y) \|f'_y(x_h, y)\|^{-1} dS + o_3(\varepsilon_k) = \frac{1}{P_h} v'(x_h) + o_3(\varepsilon_k).$$

Тогда выражение в (6) примет вид

$$(v'(x_h), Es_h) = \frac{1}{P_h} \|v'(x_h)\|^2 + o(\varepsilon_k).$$

Следовательно, существует $\delta_h > 0$ такое, что если $\varepsilon_k < \delta_h$, то $(v'(x_h), Es_h) \geq 0$. Теорема доказана.

Покажем, что значения δ_h , при которых алгоритм (5), (7), (8) является псевдоградиентным, можно вывести и явно. Она, хотя и довольно грубая оценка, определяется следующим образом.

Теорема 4. Пусть существуют предельные значения

$$M = \max_y (p(y), \|f'_y(x, y)\|, \|f'_x(x, y)\|, \|f''_{yy}(x, y)\|), \\ K = \min_y \|f'_y(x, y)\|, \quad L = \max(L_1, L_2, L_3),$$

где L_1, L_2, L_3 постоянные Липшица по y функций $p(y)$, $f'_y(x, y)$ и $f'_x(x, y)$ соответственно и пусть

$$T_h = \int \dots \int_{S(x_k, 0)} dS \left(\int \dots \int_{S(x_k, 0)} p(y) dS \right)^{-1}.$$

Тогда алгоритм (5), (7), (8) будет псевдоградиентным, если

$$\varepsilon_k \leq \delta_k = \min \left(1, \frac{K}{L_2}, \frac{1}{C_h} \|v'(x_h)\| \right), \quad (10)$$

где

$$C_h = 3T_h K^{-2} L C (2T_h K M L C + 1), \quad (11)$$

$$C = m! (1 + 12K^{-4} M^3) + 1. \quad (12)$$

Доказательство этой теоремы весьма громоздко и малоинтересно и поэтому здесь не приводится.

С вычислительной точки зрения предпочтительнее выбрать ε_k как можно большей. Чем больше ε_k , тем меньше реализаций y_k , не удовлетворяющих неравенству (8), потребуется отвергнуть. При этом уменьшится соответственно и объем вычислений. С другой стороны, из неравенства (10) следует, что для установления алгоритма (5), (7), (8) сходящимся псевдоградиентным алгоритмом необходимо выбрать $\varepsilon_k \rightarrow 0$ при $k \rightarrow \infty$.

Отметим, наконец, что требование $\delta_k \leq 1$ из условия (10) несущественно. Это требование введено лишь для упрощения оценки (10) — (12).

Автор выражает признательность И. Петерсену и Т. Тобиасу за ценные замечания.

ЛИТЕРАТУРА

1. Красносельский М. А., Забрейко П. П., Пустыльник Е. И., Соболевский П. Е., Интегральные операторы в пространствах суммируемых функций, М., 1966.

2. Robbins H., Monro S., A stochastic approximation method, *Annals Math. Stat.*, 22, No. 1, 400 (1955).
3. Kiefer E., Wolfowitz J., Stochastic estimation of the maximum of a regression function, *Annals Math. Stat.*, 23, No. 3, 462 (1952).
4. Поляк Б. Т., Цыпкин Я. З., Псевдоградиентные алгоритмы адаптации и обучения, *Автомат. и телемех.*, № 3, 45 (1973).
5. Шварц Л., *Анализ*, т. I, М., 1972.

Институт кибернетики
Академии наук Эстонской ССР

Поступила в редакцию
13/V 1974

E. RAIK

TÖENÄOSUSFUNKTSIOONI DIFERENTSEERUVUS PARAMEETRI JÄRGI JA TEMA OPTIMEERIMINE STOHHASTILISE PSEUDOGRADIENDIMEETODI ABIL

Esitatakse tõenäosusfunktsiooni diferentseerimise tingimused parameetri järgi ja gradiendi kaju ning konstrueeritakse stohhastiline pseudograadiendimeetod tõenäosusfunktsiooni minimeerimiseks.

E. RAIK

THE DIFFERENTIABILITY IN THE PARAMETER OF THE PROBABILITY FUNCTION AND OPTIMIZATION OF THE PROBABILITY FUNCTION VIA THE STOCHASTIC PSEUDOGRADIENT METHOD

The differentiability conditions of the probability function are given, the form of the gradient of this function is deduced and the stochastic pseudogradient method for optimisation presented.