

Метод Hessian Free

Авторы: Бадави Полина, Козырев Дмитрий

Оглавление

Теоретическая часть.....	4
Введение.....	4
Градиентный спуск.....	4
Метод Ньютона.....	5
Одномерный случай.....	5
Многомерный случай	5
Проблемы метода Ньютона.....	6
Метод сопряженных градиентов	6
Метод Hessian Free	8
Обход вычисления матрицы Гессе в методе Hessian Free.....	8

Теоретическая часть

Введение

Обучение нейронных сетей заключается в минимизации ошибок по отношению к наборам параметров. Как правило, большие нейронные сети могут иметь миллионы параметров, поэтому их обучение является достаточно сложной задачей.

Удивительно, но многие из последних достижений в области нейронных сетей связаны не с повышением скорости обработки данных, улучшением архитектур нейронных сетей, алгоритмов машинного обучения и прочих хитростей, а использованием мощных методов минимизации многомерной функции.

В данной работе будут рассмотрены несколько простых методов локальной минимизации многомерной функции, что позволит, в итоге, подойти к рассмотрению достаточно сложного, но эффективного метода Hessian Free.

Градиентный спуск

Простейшим итерационным алгоритмом локальной минимизации дифференцируемой функции является метод **градиентного спуска**. Градиентом функции $f(x)$ называется следующий вектор из частных производных:

$$\nabla f(x) = \left(\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right)$$

Градиент функции $f(x)$, вычисленный в определенной точке, указывает в направлении максимального роста функции в этой точке, или, что эквивалентно, в противоположном направлении к максимальному уменьшению функции $f(x)$. Формально, алгоритм запишется следующим образом:

Шаг 1. Задается начальное приближение $x^{(0)}$ и требуемая точность ε :

$$x^{(0)} = (x_1^{(0)}, x_2^{(0)}, \dots, x_n^{(0)})$$

Шаг 2. Вычисляется градиент функции в текущей точке $x^{(i)}$:

$$v^{(i)} = \nabla f(x^{(i)}).$$

Шаг 3. Осуществляется перемещение из текущей точки $x^{(i)}$ с некоторым фиксированным шагом α в направлении, противоположном вектору $v^{(i)}$:

$$x^{(i+1)} = x^{(i)} - \alpha v^{(i)}$$

Шаг 4. Проверяется условие достижения необходимой точности ε :

$$|f(x^{(i+1)}) - f(x^{(i)})| < \varepsilon$$

Если условие не выполняется, то осуществляется переход к следующей итерации на шаг 2, иначе $x^{(i+1)}$ – найденная точка минимума с точностью ε .

Существенный недостаток этого метода заключается в том, что метод градиентного спуска является методом первого порядка, то есть, в нем учитываются только первые частные производные функции $f(x)$, тогда как функция $f(x)$ может иметь достаточно сложный нелинейный вид.

Ограничение в виде игнорирования поверхностей ошибок, по своей структуре более сложных, чем плоскостей, не позволяет применять данный метод для быстрого и эффективного обучения нейронных сетей.

Метод Ньютона

Данный метод учитывает вторые производные функции $f(x)$, поэтому этот метод является методом второго порядка. Рассмотрим минимизацию методом Ньютона применительно к одномерной функции, а потом обобщим на многомерный случай.

Одномерный случай

Предположим, что функция $f(x)$ имеет минимум. Разложим функцию $f(x)$ в ряд Тейлора в точке x_0 :

$$f(x_0 + x) \approx f(x_0) + f'(x_0)x + \frac{f''(x_0)}{2}x^2 + O(x^3)$$

Мы хотим найти точку x , в которой первая производная функции в точке x равна нулю (условие экстремума функции. Используя для этого разложение функции в ряд Тейлора и отбросив члены третьего порядка и больше, получим:

$$\frac{d}{dx} \left(f(x_0) + f'(x_0)x + \frac{f''(x_0)}{2}x^2 \right) = f'(x_0) + f''(x_0)x = 0 \Rightarrow x = -\frac{f'(x_0)}{f''(x_0)}$$

Если $f(x)$ просто квадратичная функция, то это будет абсолютный минимум. Для нахождения минимума любой нелинейной функции необходимо использовать итерационный процесс, вычисляя на каждой итерации следующее приближение $x^{(i+1)}$ по формуле:

$$x^{(i+1)} = x^{(i)} - \frac{f'(x^{(i)})}{f''(x^{(i)})} = x^{(i)} - \left(f''(x^{(i)}) \right)^{-1} \left(f'(x^{(i)}) \right)$$

Если минимум существует, то алгоритм итерационно приблизится к нему с требуемой точностью.

Многомерный случай

Предположим, что для функции $f(x)$ существует минимум. В многомерном случае применяется аналогичный подход: производная эквивалентна градиенту, а вторая производная эквивалентна матрице Гессе – матрице вторых производных:

$$f'(x) \rightarrow \nabla f(x)$$

$$f''(x) \rightarrow H(f)(x)$$

где:

$$\nabla f(x) = \left(\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right)$$

$$H(f) = \begin{pmatrix} \frac{\partial^2 f}{\partial x_1 \partial x_1} & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \dots & \frac{\partial^2 f}{\partial x_n \partial x_n} \end{pmatrix}$$

Тогда формула для получения следующего приближения запишется как:

$$x^{(i+1)} = x^{(i)} - \left(H(f)(x^{(i)}) \right)^{-1} \nabla f(x^{(i)})$$

Проблемы метода Ньютона

Метод Ньютона является методом второго порядка, поэтому может работать значительно лучше, чем градиентный спуск. Предположение, что функция $f(x)$ является квадратичной, позволяет осуществлять большие шаги по направлению к минимуму при низкой кривизне (когда значение $f''(x^{(n)})$ мало), и маленькие шаги при большой кривизне. Как и в методе градиентного спуска, перемещение от текущего приближения к следующему осуществляется по направлению, противоположному направлению максимального роста функции $(-\nabla f(x^{(n)}))$.

Однако, метод Ньютона имеет очень большой недостаток: он требует вычисления матрицы Гессе, что, в свою очередь, требует $O(n^2)$ по памяти и времени, и вычисления обратной к ней, что можно реализовать методом Гаусса с прямым и обратным ходом за $O(n^3)$ по времени.

В дальнейшей работе будет показано, что метод Hessian Free лишен этих недостатков, хотя и использует основную идею метода Ньютона – приближение функции $f(x)$ ее рядом Тейлора до членов второго порядка включительно.

Метод сопряженных градиентов

Пусть $f(x)$ – квадратичная функция вида:

$$f(x) = \frac{1}{2} x^T A x + b^T x + c, \text{ где } A \in \mathbb{R}^{n \times n}, x, b, c \in \mathbb{R}^n$$

Любую квадратичную форму мы можем преобразовать к такому виду, что $A^T = A$: выпишем квадратическую форму в виде суммирования, посчитаем коэффициенты при $x_i x_j$, поделим коэффициенты на два и составим из них симметрическую матрицу A .

Градиент находится как:

$$\nabla f(x) = Ax + b$$

Пусть $x^{(0)}$ – начальное приближение. Направление, в котором находится следующее приближение, противоположно градиенту:

$$d^{(0)} = -\nabla f(x^{(0)})$$

Тогда следующее приближение вычисляется по формуле, аналогичной в методе градиентного спуска:

$$x^{(1)} = x^{(0)} + \alpha d^{(0)}$$

В отличие от обычного градиентного спуска, где шаг был фиксированным числом, в наискорейшем градиентном спуске на каждой итерации вычисляется оптимальный шаг. Вводится следующая функция, зависящая от α :

$$\begin{aligned} g(\alpha) &= f(x^{(i)} + \alpha d^{(i)}) = \frac{1}{2}(x^{(i)} + \alpha d^{(i)})^T A(x^{(i)} + \alpha d^{(i)}) + b^T(x^{(i)} + \alpha d^{(i)}) + c = \\ &= \frac{1}{2}\alpha^2 (d^{(i)})^T A d^{(i)} + (d^{(i)})^T (Ax^{(i)} + b)\alpha + \left(\frac{1}{2}(x^{(i)})^T Ax^{(i)} + (x^{(i)})^T d^{(i)} + c\right) \end{aligned}$$

Вычисление оптимального шага эквивалентно минимизации функции $g(\alpha)$. Предположим, что она имеет минимум, тогда этот минимум является глобальным (т.к. функция квадратическая). Условие на минимум:

$$g'(\alpha) = 0 = ((d^{(i)})^T A d^{(i)})\alpha + (d^{(i)})^T (Ax^{(i)} + b) \Rightarrow \alpha = -\frac{d^{(i)}(Ax^{(i)} + b)}{(d^{(i)})^T A d^{(i)}}$$

Далее необходимо учесть тот факт, чтобы следующее направление было сопряжено с текущим направлением, иначе могут возникать ситуации, в которых направления будут противоположны, поэтому мы требуем условие сопряженности.

Два вектора x и y сопряжены, если верно следующее условие: $x^T A y = 0$.

Следующее направление $d^{(1)}$, находится по формуле:

$$d^{(1)} = -\nabla f(x^{(1)}) + \beta_0 d^{(0)}$$

Здесь β_0 выбирается из условия сопряженности на вектора $d^{(1)}$ и $d^{(0)}$:

$$d^{(1)} A d^{(0)} = 0 \Rightarrow (d^{(1)})^T A d^{(0)} = -(\nabla f(x^{(1)}))^T A d^{(0)} + \beta_0 (d^{(0)})^T A d^{(0)} = 0 \Rightarrow \beta_0 = \frac{(\nabla f(x^{(1)}))^T A d^{(0)}}{(d^{(0)})^T A d^{(0)}}$$

Стоит заметить, что вычисление оптимального β_i позволяет получить направление, сопряженное со всеми предыдущими направлениями, поэтому достаточно сделать n итераций, где n – размерность пространства.

Тогда метод сопряженных градиентов для квадратичных функций формально можно записать следующим образом:

Пусть $f(x)$ – квадратичная функция вида $f(x) = \frac{1}{2}x^T A x + b^T x + c$.

Шаг 1: Инициализация. Положим $i = 0$, $x^{(i)} = x^{(0)}$ и вычислим $d^{(i)} = d^{(0)} = -\nabla f(x^{(0)})$.

Шаг 2: Поиск оптимального шага. Вычислим α как результат минимизации функции $f(x^{(i)} + \alpha d^{(i)})$ по формуле:

$$\alpha = -\frac{(d^{(i)})^T (Ax^{(i)} + b)}{(d^{(i)})^T A d^{(i)}}$$

Шаг 3. Обновление приближения. Вычислим $x^{(i+1)} = x^{(i)} + \alpha d^{(i)}$

Шаг 4. Обновление направления. Вычислим $d^{(i+1)} = -\nabla f(x^{(i+1)}) + \beta_i d^{(i)}$, где β_i находится как:

$$\beta_i = \frac{(\nabla f(x^{(i+1)}))^T Ad^{(i)}}{(d^{(i)})^T Ad^{(i)}}$$

Шаг 5. Итерационный процесс. Повторим шаги 2-4 пока мы не рассмотрим n направлений d_i , где n – размерность пространства.

Метод Hessian Free

Алгоритмически метод можно записать в следующей форме:

Пусть задана функция $f(x): \mathbb{R}^n \rightarrow \mathbb{R}$, которую необходимо минимизировать, ε – необходимая точность, и точка $x^{(0)}$ – начальное приближение.

Шаг 1. Инициализация. Положим $i = 0$ и $x^{(i)} = x^{(0)}$.

Шаг 2. Аппроксимация квадратичной функцией. Для текущего приближения $x^{(i)}$ вычислим градиент $\nabla f(x^{(i)})$ и матрицу Гессе $H(f)(x^{(i)})$, и будем иметь в виду, что для функции $f(x)$ справедливо разложение в ряд Тейлора:

$$f(x + \Delta x) \approx f(x) + (\nabla f(x))^T \Delta x + (\Delta x)^T H(f) \Delta x + O(\|\Delta x\|^3)$$

Шаг 3. Сопряженный градиент. Вычислим $x^{(i+1)}$ используя метод сопряженных градиентов для квадратичной функции (для текущего разложения в ряд Тейлора). В качестве переменной для метода сопряженных градиентов служит Δx .

Шаг 4. Итерационный процесс. Повторим шаги 2 и 3 пока $x^{(i)}$ пока процесс не сойдется с необходимой точностью ε .

Данный алгоритм вбирает в себя все описанные ранее идеи. Рассмотрим обход вычисления матрицы Гесса на шаге 2.

Обход вычисления матрицы Гессе в методе Hessian Free

В методе Ньютона была необходимость в вычислении матрицы Гессе, но в данном алгоритме вычисление матрицы Гессе не требуется, так как во всех формулах используется результат умножения матрицы Гессе на вектор.

Рассмотрим i – ый элемент вектора, получающегося после произведения матрицы Гессе на вектор:

$$(H(f)v)_i = \sum_{j=1}^N \frac{\partial^2 f}{\partial x_i \partial x_j}(x) \cdot v_j = \nabla \frac{\partial f}{\partial x_i}(x) \cdot v$$

Следует заметить, что это ни что иное как производная по направлению от $\frac{\partial f}{\partial x_i}$ в направлении v :

$$\nabla_v f = \lim_{\varepsilon \rightarrow 0} \frac{f(x + \varepsilon v) - f(x)}{\varepsilon}$$

Используя разностную аппроксимацию дифференцирования для малого шага h , получим приближение с точностью $O(\|hv\|)$ или $O(\|hv\|^2)$:

$$\nabla_v f \approx \frac{f(x + hv) - f(x)}{h} + O(|h| \|v\|^2)$$

$$\nabla_v f \approx \frac{f(x + hv) - f(x - hv)}{2h} + O(|h|^2 \|v\|^3)$$

При желании, можно доказать эти формулы, используя разложение Тейлора функции $f(x)$.

Тогда аппроксимацию умножения матрицы Гессе на вектор можно вычислить как:

$$Hv \approx \frac{\nabla f(x + hv) - \nabla f(x)}{h} + O(|h| \|v\|^2)$$

$$Hv \approx \frac{\nabla f(x + hv) - \nabla f(x - hv)}{2h} + O(|h|^2 \|v\|^3)$$