

Градиентный спуск

Нейронные сети

Теория игр, 2023



- 1 Машинное обучение с учителем
- 2 Градиентный спуск
- 3 Искусственная нейронная сеть



1 Машинное обучение с учителем

2 Градиентный спуск

3 Искусственная нейронная сеть



- Christopher Bishop. Pattern Recognition and Machine Learning
- Shai Shalev-Shwartz, Shai Ben-David. Understanding Machine Learning: From Theory to Algorithms
- Jorge Nocedal, Stephen J. Wright. Numerical Optimization



Множество объектов $\mathcal{X} = \{\mathbf{x}\}$, представленных вектором из d признаков $\mathbf{x} = (x_1, \dots, x_d)$

Множество значений зависимой переменной $\mathcal{Y} = \{y\}$

Множество моделей(гипотез) \mathcal{H} , описываемых в виде функции, которая каждому объекту ставит в соответствие одно значений зависимой переменной $h : \mathcal{X} \rightarrow \mathcal{Y}$.

Функция потерь $\ell : \mathcal{H} \times \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}_+$



Истинная ошибка

Множество объектов $\mathcal{X} = \{\mathbf{x}\}$, представленных вектором из d признаков $\mathbf{x} = (x_1, \dots, x_d)$

Множество значений зависимой переменной $\mathcal{Y} = \{y\}$

Множество моделей (гипотез) \mathcal{H} , описываемых в виде функции, которая каждому объекту ставит в соответствие одно значений зависимой переменной $h : \mathcal{X} \rightarrow \mathcal{Y}$.

Функция потерь $\ell : \mathcal{H} \times \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}_+$

При заданном **вероятностном распределении** \mathcal{D} над $\mathcal{X} \times \mathcal{Y}$ функция ℓ является случайной

Истинная ошибка модели $h \in \mathcal{H}$

математическое ожидание функции потери модели h :

$$L_D(h) = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\ell(h(\mathbf{x}), y)]$$



Эмпирическая ошибка

Множество объектов $\mathcal{X} = \{x\}$, представленных вектором из d признаков $x = (x_1, \dots, x_d)$

Множество значений зависимой переменной $\mathcal{Y} = \{y\}$

Множество моделей (гипотез) \mathcal{H} , описываемых в виде функции, которая каждому объекту ставит в соответствие одно значений зависимой переменной $h : \mathcal{X} \rightarrow \mathcal{Y}$.

Функция потерь $\ell : \mathcal{H} \times \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}_+$

Так как в большинстве случаев распределение \mathcal{D} неизвестно, то модель выбирается по **имеющейся выборке из m объектов**

$$S = ((x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)}))$$

Эмпирическая ошибка модели $h \in \mathcal{H}$

$$L_S(h) = \frac{1}{m} \sum_{i=1}^m \ell(h(x^{(i)}), y^{(i)}) \quad (1)$$

Задача регрессии

Множество объектов $\mathcal{X} = \{\mathbf{x}\}$, представленных вектором из d признаков $\mathbf{x} = (x_1, \dots, x_d)$

Множество значений зависимой переменной $\mathcal{Y} = \mathbb{R}$

Множество моделей(гипотез) \mathcal{H} , описываемых в виде функции, которая каждому объекту ставит в соответствие одно значений зависимой переменной $h : \mathcal{X} \rightarrow \mathcal{Y}$.

Функция потерь $\ell : \mathcal{H} \times \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}_+$

$$h, \mathbf{x}, y \mapsto (h(\mathbf{x}) - y)^2$$

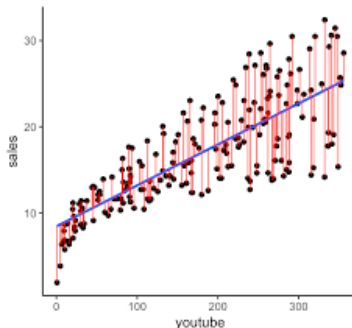
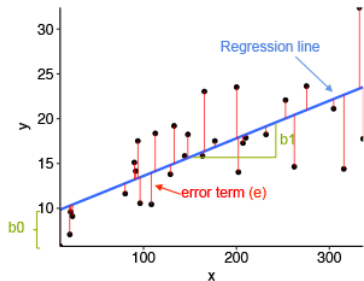


Пример задачи регрессии

Объект задается набором (вектором) d признаков $\mathbf{x} = (x_1, \dots, x_d)$

Модель - линейная функция от вектора параметров \mathbf{w}

$$h_{\mathbf{w}}(\mathbf{x}) = h(\mathbf{x}, \mathbf{w}) = w_0 + x_1 * w_1 + \dots + x_d * w_d = \langle \mathbf{x}, \mathbf{w} \rangle$$



Задача классификации

Множество объектов $\mathcal{X} = \{x\}$, представленных вектором из d признаков $x = (x_1, \dots, x_d)$

Множество значений зависимой переменной $\mathcal{Y} = y_1, \dots, y_n$

Множество моделей(гипотез) \mathcal{H} , описываемых в виде функции, которая каждому объекту ставит в соответствие одно значений зависимой переменной $h : \mathcal{X} \rightarrow \mathcal{Y}$.

Функция потерь $\ell : \mathcal{H} \times \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}_+$

$$\ell(h(x), y) \mapsto \begin{cases} 1, & \text{если } h(x) \neq y \\ 0, & \text{иначе} \end{cases}$$



Пример задачи классификации

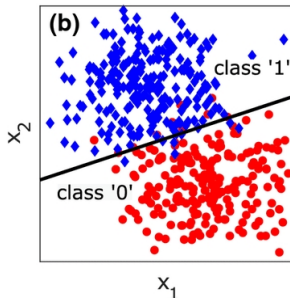
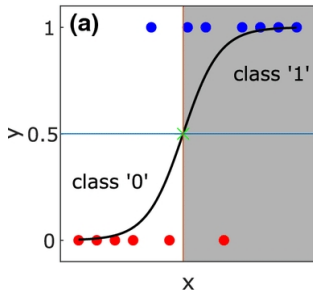
Объект задается набором (вектором) d признаков $\mathbf{x} = (x_1, \dots, x_d)$

Множество значений зависимой переменной - два класса $\mathcal{Y} = \{0, 1\}$

Модель - логистическая функция от вектора параметров \mathbf{w}

$$h_{\mathbf{w}}(\mathbf{x}) = h(\mathbf{x}, \mathbf{w}) = \frac{1}{1 + e^{\langle \mathbf{x}, \mathbf{w} \rangle}}, \langle \mathbf{x}, \mathbf{w} \rangle = w_0 + x_1 * w_1 + \dots + x_d * w_d$$

$$y_{\text{гипотеза}} = \begin{cases} 1, & \text{если } h_{\mathbf{w}}(\mathbf{x}) > 0.5 \\ 0, & \text{иначе} \end{cases}$$



- 1 Машинное обучение с учителем
- 2 Градиентный спуск**
- 3 Искусственная нейронная сеть



Минимизация эмпирической ошибки модели $h \in \mathcal{H}$

$$L_S(h_{\mathbf{w}}) = L_S(\mathbf{w}) = \frac{1}{m} \sum_{i=1}^m \ell(h(\mathbf{x}^{(i)}, \mathbf{w}), y^{(i)}) \rightarrow \min_{\mathbf{w}} \quad (2)$$

Общая постановка задачи оптимизации для произвольной функции $f : \mathbb{R}^n \rightarrow \mathbb{R}$

$$f(\mathbf{w}) \rightarrow \min_{\mathbf{w}} \quad (3)$$

\mathbf{w}^* - минимум функции f , если $\forall \mathbf{w} \ f(\mathbf{w}^*) \leq f(\mathbf{w})$



$$f : \mathbb{R}^n \rightarrow \mathbb{R}$$

$$\mathbf{w} \mapsto f(\mathbf{w})$$

Градиент функции f в точке p

$$\nabla f : \mathbb{R}^n \rightarrow \mathbb{R}^n$$

$$\nabla f_{\mathbf{w}}(p) = \begin{bmatrix} \frac{\partial f}{\partial w_1}(p) \\ \frac{\partial f}{\partial w_2}(p) \\ \vdots \\ \frac{\partial f}{\partial w_n}(p) \end{bmatrix}$$



Нахождение локального минимума функции $f : \mathbb{R}^n \rightarrow \mathbb{R}$

Алгоритм

Инициализация \mathbf{w}^0 произвольно

Для каждой итерации

выбор шага обновления α

$$\mathbf{w}^{k+1} \leftarrow \mathbf{w}_k - \alpha \nabla f(\mathbf{w}_k)$$

Для минимума \mathbf{w}^* верно $\nabla f(\mathbf{w}^*) = 0$

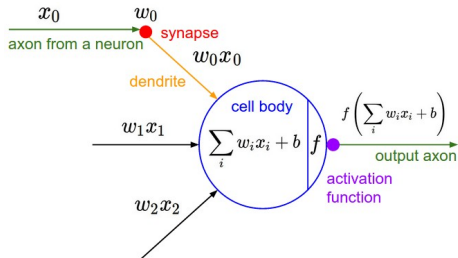
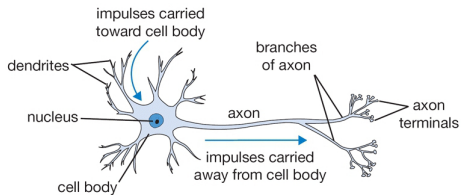


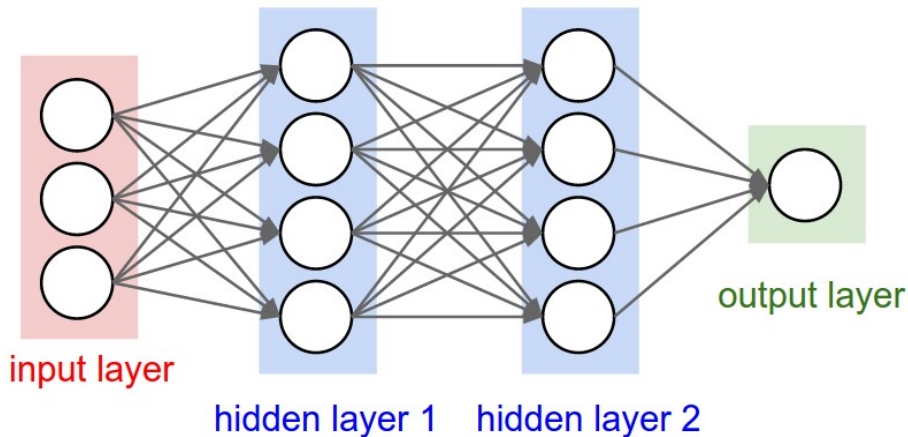
- 1 Машинное обучение с учителем
- 2 Градиентный спуск
- 3 Искусственная нейронная сеть**



Нейрон

<https://cs231n.github.io>



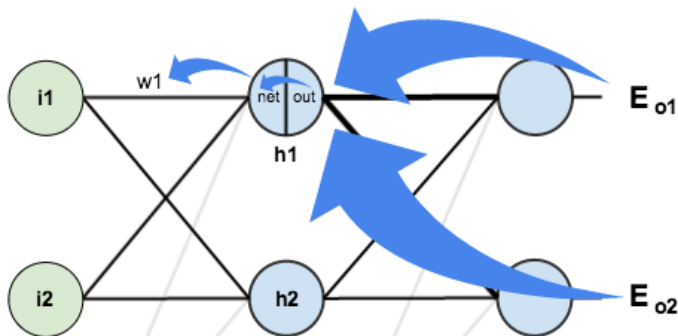


Алгоритм обратного распространения

$$\frac{\partial E_{total}}{\partial w_1} = \frac{\partial E_{total}}{\partial out_{h1}} * \frac{\partial out_{h1}}{\partial net_{h1}} * \frac{\partial net_{h1}}{\partial w_1}$$

↓

$$\frac{\partial E_{total}}{\partial out_{h1}} = \frac{\partial E_{o1}}{\partial out_{h1}} + \frac{\partial E_{o2}}{\partial out_{h1}}$$



$$E_{total} = E_{o1} + E_{o2}$$