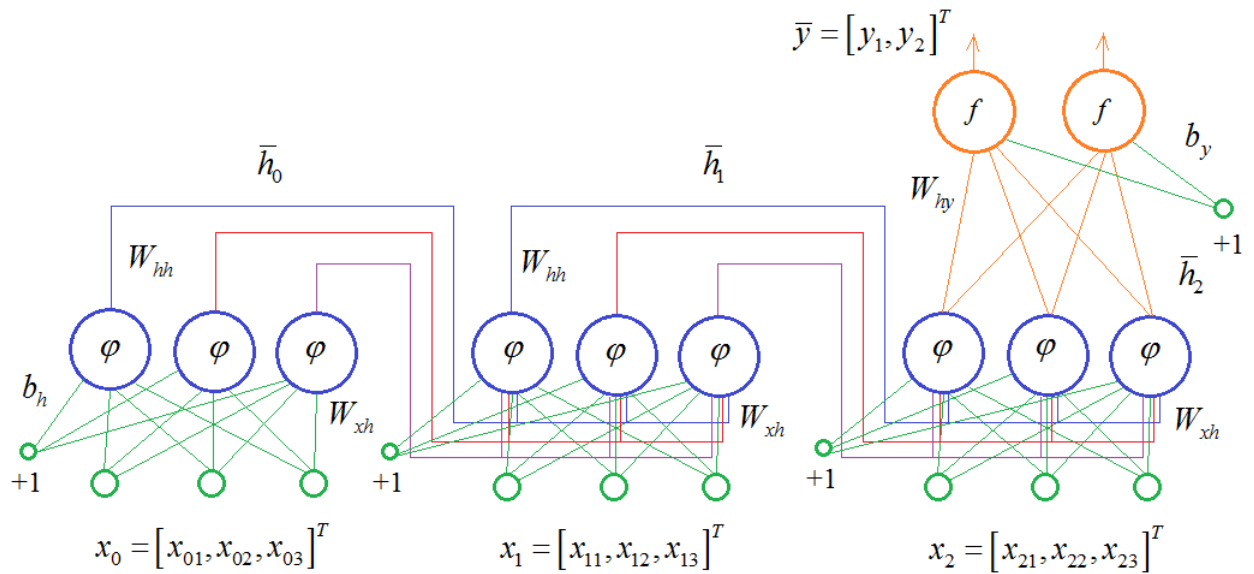


## Пример алгоритма обратного распространения ошибки по времени (Back Propagation Through Time)



Развернем рекуррентную НС на три шага во времени. В каждый момент времени на ее вход подается вектор  $\bar{x}$ , и на третьем шаге смотрим выходное значение  $\bar{y}$ . Для обучения такой сети можно использовать алгоритм

back propagation

с учетом временного характера поведения сети. В нашем случае рекуррентная сеть построена по принципу

many to one (многие к одному)

то есть, множество входных сигналов и один выходной. Так как мы рассматриваем задачу классификации, то в качестве функции активации выходных нейронов выберем softmax, а потери будем считать через кросс-энтропию:

$$E = -\sum_{i=1}^M t_i \log(y_i) = -\sum_{i=1}^M t_i \log(\text{softmax}(v_i))$$

где  $M = 2$  – число выходных нейронов;  $v_i$  - индуцированный (суммарный) сигнал на входах нейронов последнего слоя;  $t_i = \{0;1\}$  - требуемые выходные значения (в нашем случае 1 или 0);  $\log$  – логарифм, например, натуральный.

Первым делом нам нужно вычислить градиент функции потерь  $E$  от входного значения  $v_i$ . Для строго определенного класса

$$i = C$$

функция

$$E = -\log(\text{softmax}(v_C))$$

Следовательно:

$$\frac{\partial E}{\partial v_i} = \begin{cases} y_i - 1, & i = C \\ y_i, & i \neq C \end{cases}$$

Например, если на выходе наблюдаем значения:

$$\bar{y} = [0,8; 0,2]^T$$

то  $\frac{\partial E}{\partial \bar{v}}$  при  $C = 1$  образует вектор:

$$\frac{\partial E}{\partial \bar{v}} = [0,8 - 1; 0,2]^T = [-0,2; 0,2]^T$$

Как видите, все предельно просто. Далее, определим градиент для весов матрицы  $W_{hy}$ :

$$\frac{\partial E}{\partial W_{hy}} = \frac{\partial E}{\partial \bar{v}} \cdot \frac{\partial \bar{v}}{\partial W_{hy}}$$

Так как

$$\bar{v} = W_{hy} \cdot \bar{h}_n + b_y$$

(здесь  $\bar{h}_n = \bar{h}_2$ ), то

$$\frac{\partial \bar{v}}{\partial W_{hy}} = \bar{h}_n^T$$

и градиент изменения весов:

$$\frac{\partial E}{\partial W_{hy}} = \frac{\partial E}{\partial \bar{v}} \cdot \bar{h}_n^T$$

Аналогично вычисляется градиент изменения биаса:

$$\frac{\partial E}{\partial b_y} = \frac{\partial E}{\partial \bar{v}} \cdot \frac{\partial \bar{v}}{\partial b_y} = \frac{\partial E}{\partial \bar{v}}$$

Используя эти величины, мы теперь можем применить алгоритм градиентного спуска для корректировки весов  $W_{hy}$  и смещения  $b_y$ :

$$W_{hy} = W_{hy} - \lambda \cdot \frac{\partial E}{\partial \bar{v}} \cdot \bar{h}_n^T$$

$$b_y = b_y - \lambda \cdot \frac{\partial E}{\partial \bar{v}}$$

Осталось выполнить похожие вычисления для весов матриц  $W_{hh}, W_{xh}$  и биаса  $b_h$ . Начнем с матрицы  $W_{hh}$ :

$$\frac{\partial E}{\partial W_{hh}} = \frac{\partial E}{\partial \bar{v}} \cdot \frac{\partial \bar{v}}{\partial \bar{h}_n} \cdot \frac{\partial \bar{h}_n}{\partial W_{hh}}$$

Это выражение было бы верно для обычных сетей прямого распространения. Но у нас здесь рекуррентная сеть и значение вектора:

$$\bar{h}_n = \varphi(W_{hh} \cdot \bar{h}_{n-1} + W_{xh} \cdot \bar{x}_n + b_h)$$

зависит от значения на предыдущем шаге. Поэтому для корректного вычисления градиента нам нужно пройти по всей рекурсии до самого начала:

$$\frac{\partial E}{\partial W_{hh}} = \frac{\partial E}{\partial \bar{v}} \cdot \left[ \frac{\partial \bar{v}}{\partial \bar{h}_n} \cdot \frac{\partial \bar{h}_n}{\partial W_{hh}} + \frac{\partial \bar{v}}{\partial \bar{h}_{n-1}} \cdot \frac{\partial \bar{h}_{n-1}}{\partial W_{hh}} + \dots + \frac{\partial \bar{v}}{\partial \bar{h}_0} \cdot \frac{\partial \bar{h}_0}{\partial W_{hh}} \right]$$

В результате, получаем такую формулу:

$$\frac{\partial E}{\partial W_{hh}} = \frac{\partial E}{\partial \bar{v}} \cdot \sum_{t=0}^n \frac{\partial \bar{v}}{\partial \bar{h}_t} \cdot \frac{\partial \bar{h}_t}{\partial W_{hh}}$$

Так как

$$\bar{h}_t = \varphi(W_{hh} \cdot \bar{h}_{t-1} + W_{xh} \cdot \bar{x}_t + b_t)$$

то

$$\frac{\partial \bar{h}_t}{\partial W_{hh}} = \varphi'(x) \cdot \bar{h}_{t-1}^T$$

Давайте для определенности положим, что функция активации — это гиперболический тангенс:

$$\varphi(x) = \tanh(x)$$

производная которой, равна:

$$\varphi'(x) = 1 - \tanh^2(x)$$

Тогда

$$\frac{\partial \bar{h}_t}{\partial W_{hh}} = (1 - \bar{h}_t^2) \cdot \bar{h}_{t-1}^T$$

Далее, нужно вычислить производную  $\frac{\partial \bar{v}}{\partial \bar{h}_t}$ . Так как она зависит от всех предыдущих выходов  $\bar{h}_t$ , то проще всего ее вычислять по рекурсии:

$$\frac{\partial \bar{v}}{\partial \bar{h}_t} = \frac{\partial \bar{v}}{\partial \bar{h}_{t+1}} \cdot \frac{\partial \bar{h}_{t+1}}{\partial \bar{h}_t}$$

где

$$\frac{\partial \bar{h}_{t+1}}{\partial \bar{h}_t} = \varphi'(x) \cdot W_{hh}^T = [1 - \bar{h}_t^2] \cdot W_{hh}^T$$

а

$$\frac{\partial \bar{v}}{\partial h_n} = W_{hy}^T$$

Имея эти выражения, мы теперь можем вычислить суммарный градиент по всей рекурсии (то есть по времени) для матрицы весов  $W_{hh}$ :

$$\frac{\partial E}{\partial W_{hh}} = \frac{\partial E}{\partial \bar{v}} \cdot \sum_{t=0}^n \frac{\partial \bar{v}}{\partial \bar{h}_t} \cdot (1 - \bar{h}_t^2) \cdot \bar{h}_{t-1}^T$$

Аналогично вычисляются следующие градиенты:

$$\frac{\partial E}{\partial W_{xh}} = \frac{\partial E}{\partial \bar{v}} \cdot \sum_{t=0}^n \frac{\partial \bar{v}}{\partial \bar{h}_t} \cdot (1 - \bar{h}_t^2) \cdot \bar{x}_t^T$$

$$\frac{\partial E}{\partial b_h} = \frac{\partial E}{\partial \bar{v}} \cdot \sum_{t=0}^n \frac{\partial \bar{v}}{\partial \bar{h}_t} \cdot (1 - \bar{h}_t^2)$$

Используя эти выражения, выполняем корректировку соответствующих весов:

$$W_{hh} = W_{hh} - \lambda \cdot \frac{\partial E}{\partial W_{hh}}$$

$$W_{xh} = W_{xh} - \lambda \cdot \frac{\partial E}{\partial W_{xh}}$$

$$b_h = b_h - \lambda \cdot \frac{\partial E}{\partial b_h}$$

Вот так выглядит алгоритм обратного распространения ошибки по времени для рекуррентных нейронных сетей.