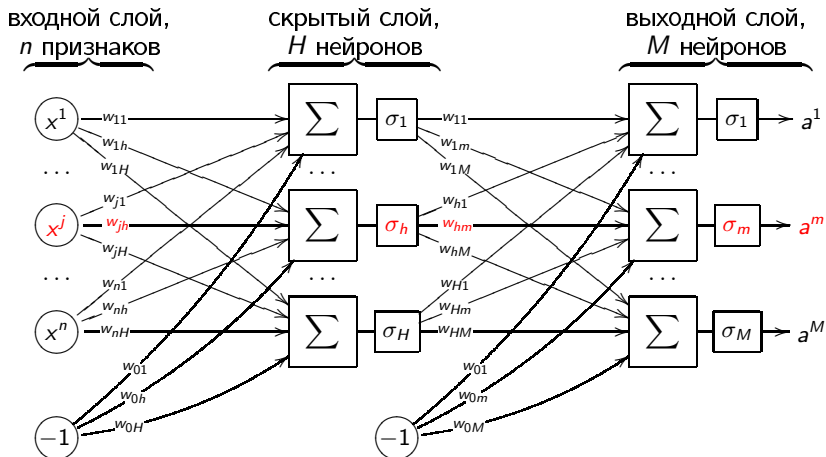


# Многослойная нейронная сеть

Пусть для общности  $Y = \mathbb{R}^M$ , для простоты слоёв только два:

$$a^m(x, w) = \sigma_m \left( \sum_{h=0}^H w_{hm} \sigma_h \left( \sum_{j=0}^n w_{jh} f_j(x_i) \right) \right).$$



Задача минимизации суммарных потерь:

$$Q(w) := \sum_{i=1}^{\ell} \mathcal{L}_i(w) \rightarrow \min_w .$$

**Вход:** выборка  $X^\ell$ ; темп обучения  $\eta$ ; параметр  $\lambda$ ;

**Выход:** веса  $w \equiv (w_{jh}, w_{hm}) \in \mathbb{R}^{H(n+M+1)+M}$ ;

- 1: инициализировать веса  $w$  и текущую оценку  $Q(w)$ ;
- 2: **повторять**
- 3: выбрать объект  $x_i$  из  $X^\ell$  (например, случайно);
- 4: вычислить потерю  $\mathcal{L}_i := \mathcal{L}_i(w)$ ;
- 5: градиентный шаг:  $w := w - \eta \nabla \mathcal{L}_i(w)$ ;
- 6: оценить значение функционала:  $Q := (1 - \lambda)Q + \lambda \mathcal{L}_i$ ;
- 7: **пока** значение  $Q$  и/или веса  $w$  не стабилизируются;

## Задача дифференцирования суперпозиции функций

Выходные значения сети  $a^m(x_i)$ ,  $m = 1, \dots, M$  на объекте  $x_i$ :

$$a^m(x_i) = \sigma_m \left( \sum_{h=0}^H w_{hm} u^h(x_i) \right); \quad u^h(x_i) = \sigma_h \left( \sum_{j=0}^J w_{jh} f_j(x_i) \right).$$

Пусть для конкретности  $\mathcal{L}_i(w)$  — средний квадрат ошибки:

$$\mathcal{L}_i(w) = \frac{1}{2} \sum_{m=1}^M (a^m(x_i) - y_i^m)^2.$$

**Промежуточная задача:** найти частные производные

$$\frac{\partial \mathcal{L}_i(w)}{\partial a^m}; \quad \frac{\partial \mathcal{L}_i(w)}{\partial u^h}.$$

**Промежуточная задача:** частные производные

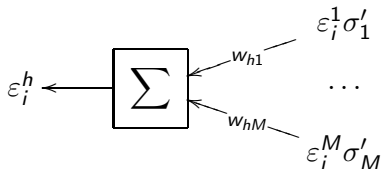
$$\frac{\partial \mathcal{L}_i(w)}{\partial a^m} = a^m(x_i) - y_i^m = \varepsilon_i^m$$

— это ошибка на выходном слое;

$$\frac{\partial \mathcal{L}_i(w)}{\partial u^h} = \sum_{m=1}^M (a^m(x_i) - y_i^m) \sigma'_m w_{hm} = \sum_{m=1}^M \varepsilon_i^m \sigma'_m w_{hm} = \varepsilon_i^h$$

— назовём это *ошибкой на скрытом слое*.

Похоже, что  $\varepsilon_i^h$  вычисляется по  $\varepsilon_i^m$  путём его пропускания через сеть в обратном направлении:



Теперь, имея частные производные  $\mathcal{L}_i(w)$  по  $a^m$  и  $u^h$ , легко выписать градиент  $\mathcal{L}_i(w)$  по весам  $w$ :

$$\frac{\partial \mathcal{L}_i(w)}{\partial w_{hm}} = \frac{\partial \mathcal{L}_i(w)}{\partial a^m} \frac{\partial a^m}{\partial w_{hm}} = \varepsilon_i^m \sigma'_m u^h(x_i), \quad m = 1..M, \quad h = 0..H;$$

$$\frac{\partial \mathcal{L}_i(w)}{\partial w_{jh}} = \frac{\partial \mathcal{L}_i(w)}{\partial u^h} \frac{\partial u^h}{\partial w_{jh}} = \varepsilon_i^h \sigma'_h f_j(x_i), \quad h = 1..H, \quad j = 0..n;$$

**Алгоритм обратного распространения ошибки BackProp:**

**Вход:**  $X^\ell = (x_i, y_i)_{i=1}^\ell \subset \mathbb{R}^n \times \mathbb{R}^M$ ; параметры  $H, \lambda, \eta$ ;

**Выход:** синаптические веса  $w_{jh}, w_{hm}$ ;

---

1: ...

- 1: инициализировать веса  $w_{jh}$ ,  $w_{hm}$ ;
- 2: **повторять**
- 3: выбрать объект  $x_i$  из  $X^\ell$  (например, случайно);
- 4: прямой ход:  
$$u_i^h := \sigma_h\left(\sum_{j=0}^J w_{jh}x_i^j\right), \quad h = 1..H;$$
$$a_i^m := \sigma_m\left(\sum_{h=0}^H w_{hm}u_i^h\right), \quad \varepsilon_i^m := a_i^m - y_i^m, \quad m = 1..M;$$
$$\mathcal{L}_i := \sum_{m=1}^M (\varepsilon_i^m)^2;$$
- 5: обратный ход:  $\varepsilon_i^h := \sum_{m=1}^M \varepsilon_i^m \sigma'_m w_{hm}$ ,  $h = 1..H$ ;
- 6: градиентный шаг:  
$$w_{hm} := w_{hm} - \eta \varepsilon_i^m \sigma'_m u_i^h, \quad h = 0..H, \quad m = 1..M;$$
$$w_{jh} := w_{jh} - \eta \varepsilon_i^h \sigma'_h x_i^j, \quad j = 0..n, \quad h = 1..H;$$
- 7:  $Q := (1 - \lambda)Q + \lambda \mathcal{L}_i$ ;
- 8: **пока**  $Q$  не стабилизируется;

## Преимущества:

- быстрое вычисление градиента;
- метод легко обобщается на любые  $\sigma$ ,  $\mathcal{L}$ ;
- возможно динамическое (потокковое) обучение;
- на сверхбольших выборках не обязательно брать все  $x_i$ ;
- возможность распараллеливания;

## Недостатки — все те же, свойственные SG:

- возможна медленная сходимость;
- застревание в локальных минимумах;
- проблема «паралича сети» (горизонтальные асимптоты  $\sigma$ );
- проблема переобучения;
- подбор комплекса эвристик является искусством;