## Связь дисперсий последовательных слоёв глубокой нейросети

Educated guess: зависит от (начальной инициализации) весов

Рассмотрим кейс:

$w_i \sim U\left[-\frac{1}{\sqrt{n_{out}}}, \frac{1}{\sqrt{n_{out}}}\right]$, где

$w_i \in W_j$ — матрица весов

$n_{out}$ — кол-во нейронов на выходном слое

$$E[w_i] = 0, \quad Var[w_i] = \frac{1}{3 n_{out}}$$

Напоминание:

$x \sim U[a, в]$ — равномерн. распред-е сл. в.

$E[x] = \frac{a+в}{2}$

$Var(x) = \frac{(в-a)^2}{12}$

Из лекции мы помним, что

$$Var\left(y_i^{(\ell+1)}\right) = Var\left(w_i^{(\ell+1)}\right) Var\left(y_i^{(\ell)}\right)$$

где $\ell$ — номер слоя

т.е. $= Var\left(w_i^{(\ell+1)}\right) Var\left(w_i^{(\ell)}\right) Var\left(y_i^{(\ell-1)}\right)$

$= ...$

$$= \prod_\ell Var\left(w_i^\ell\right) Var\left(y_i^{(0)}\right)$$

на каждом слое веса распределены одинаково и не зависят др. от др.

$$= \prod_\ell \frac{1}{3 n_{out}^{(\ell)}} \cdot Var\left(y_i^{(0)}\right) \xrightarrow[\ell \text{ большая}]{} 0$$

То есть с увеличением глубины сети сигнал затухает: нулевые матожидание и дисперсия.

Это уже достаточно плохо, но дальше только хуже. Рассмотрим обратный проход.

Посчитаем градиент функции ошибки $L$ на слое $\ell$. Вспомним, что слой $\ell+1$ — это

$$\vdots = f\left(\vdots \quad W^{(\ell+1)} - \text{веса} \quad y^{(\ell)}\right) \circ, \quad y_i^{(\ell+1)} = f\left(\sum_j w_{ij}^{(\ell+1)} y_j^{(\ell)}\right)$$

$y^{(\ell+1)}$ — $f$-ция активации

$$\frac{\partial L}{\partial y_i^{(\ell)}} = \sum_j \frac{\partial L}{\partial y_j^{(\ell+1)}} \cdot \frac{\partial y_j^{(\ell+1)}}{\partial y_i^{(\ell)}} = \sum_j \frac{\partial L}{\partial y_j^{(\ell+1)}} \boxed{\left(\frac{\partial f}{\partial y_i^{(\ell)}}\right)} \frac{\partial \left(\sum w_{ij} y_j^{(\ell)}\right)}{\partial y_i} =$$

$\uparrow$ не зависит от $j$

$$= \frac{\partial f}{\partial y_i^{(\ell)}} \sum_j \frac{\partial L}{\partial y_j^{(\ell+1)}} \cdot w_{ij}^{(\ell+1)}$$

Для симметричной относительно начала координат ф-ции активации $f$ (напр., $\tanh$)

$\dfrac{\partial f}{\partial y_i^{(\ell)}} \approx 1$ в окрестности нуля

(рассм. именно эту область, потому что у нас нулевое мат.ожидание )

$$\text{Var}\left(\frac{\partial L}{\partial y_i^{(\ell)}}\right) = \text{Var}\left(\sum_j \frac{\partial L}{\partial y_j^{(\ell+1)}} w_{ij}^{(\ell+1)}\right) = \sum_j \text{Var}\left(\frac{\partial L}{\partial y_j^{(\ell+1)}}\right) \cdot \text{Var}\, w_{ij}^{(\ell+1)}$$

$$\sim \frac{1}{3n_{in}} \cdot \sum_j \cdot \sum \text{Var}\frac{\partial L}{\partial y_j^{(\ell+1)}}$$

$\uparrow$ число входящих нейронов

И здесь дисперсия стремится к нулю.

Вывод: такая инициализация весов $\omega_i \sim u\left[-\frac{1}{\sqrt{n_{out}}}, \frac{1}{\sqrt{n_{out}}}\right]$ не подходит для глубоких сетей.

Решение

o Xavier initialization для симметричной ф-ции активации

o He initialization для ReLU.