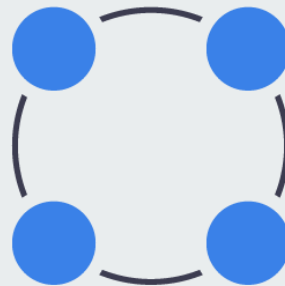




Машинное обучение

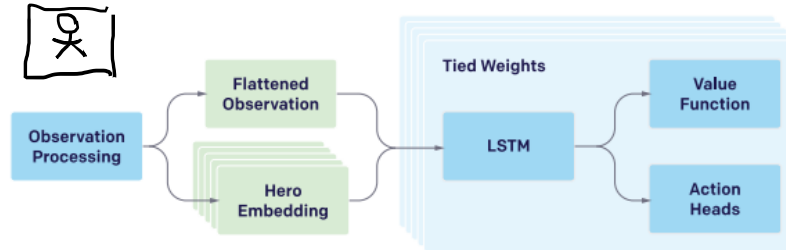
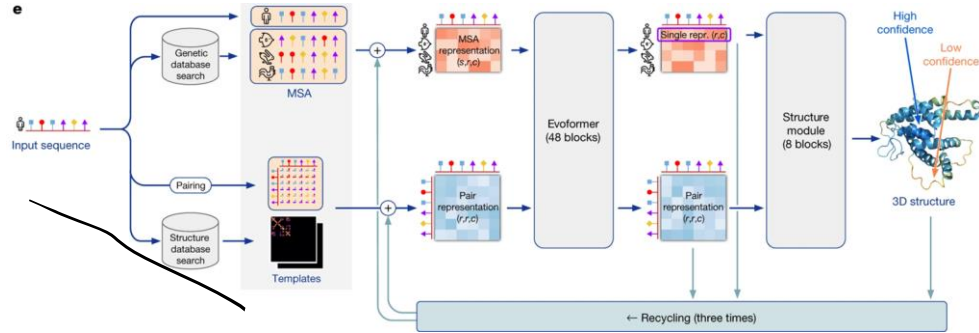
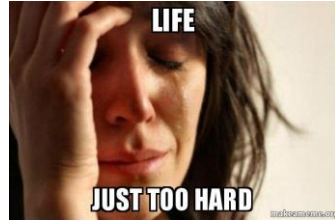
Лекция 7. Нейронные сети. Начало

(19.03.2022)

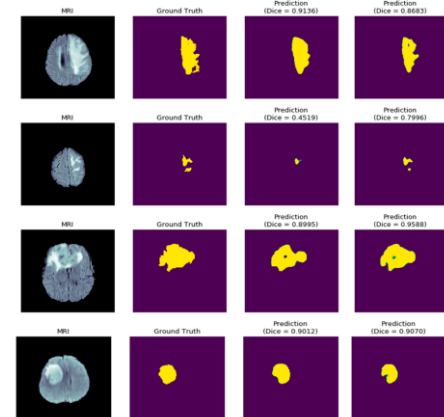
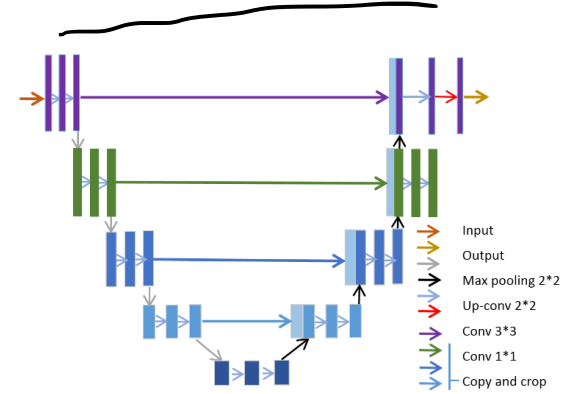




Начало.



UNET



Inventing magic spells



I trained a neural network twice on [Dungeons and Dragons spells](#), and once on [spells from Harry Potter](#). See if you can figure out which list is which.

Chorus of the dave
Song of the doom goom
Barking Sphere
Gland Growth
Hold Mouse

Hurder-gerping Charm
Regrowing hair to curse of the Bogies
Brechaim hedbivicus Doobers Spell
Fubbledory Charm
Squggly-wing fart



With all my
sparklepants!



Hugs for your
Valentine,
from the
inside!



Supermacroo-
textacularly
big thanks
for you!



Hacks, kisses
and nuzzle
nuzzles



Valentine:
How cuddly!



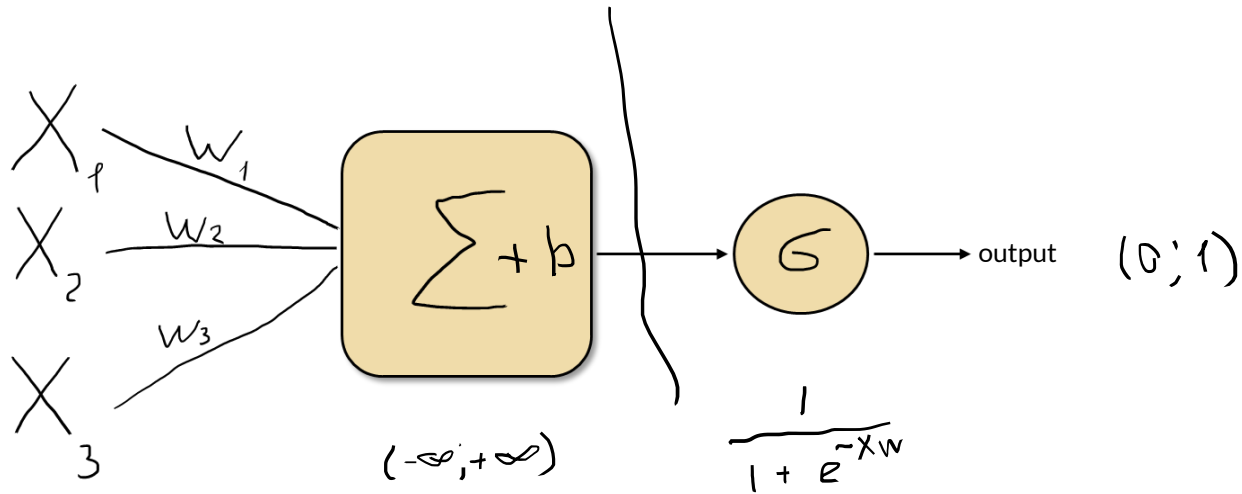
You're the
snail's
poise!



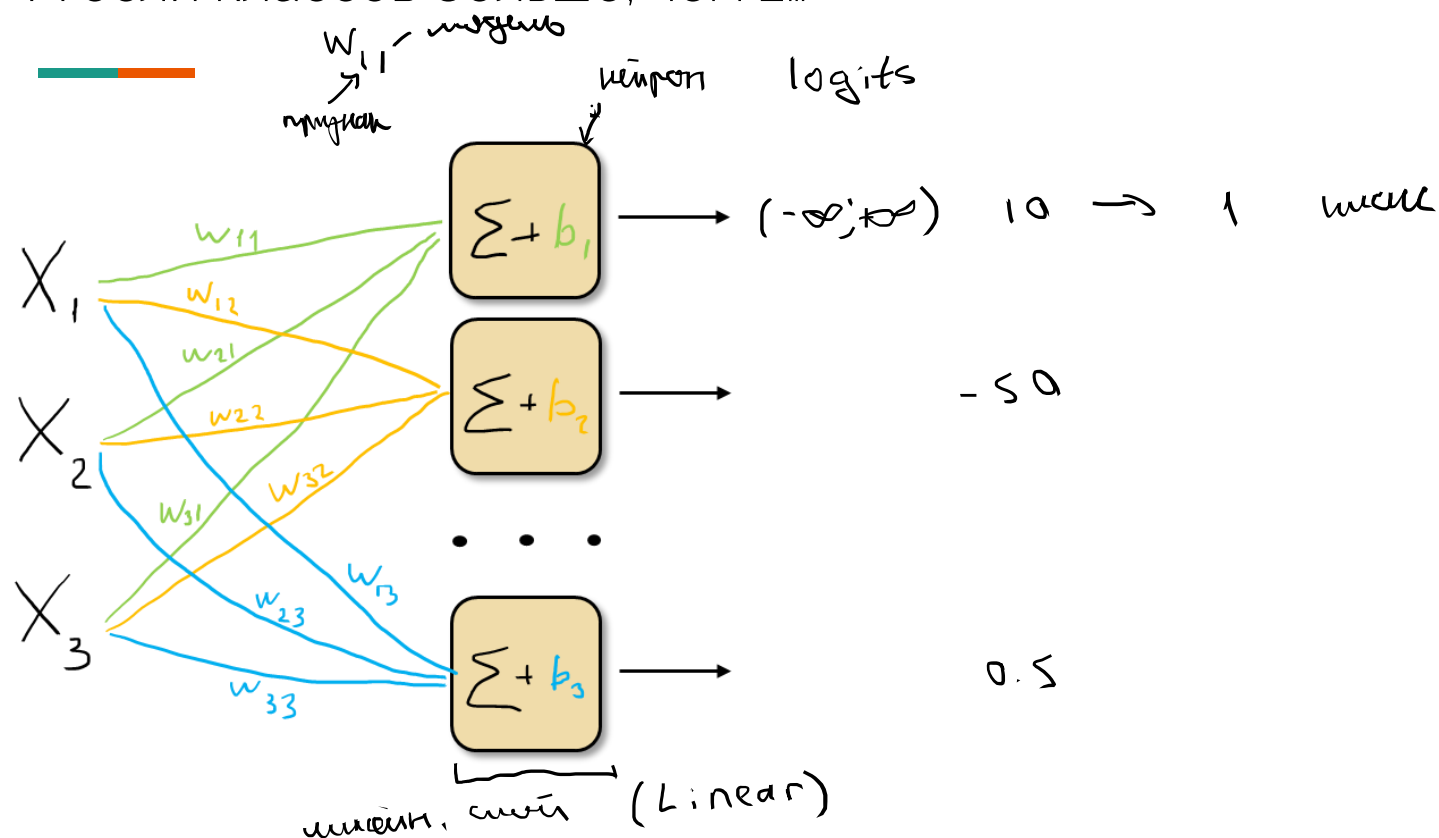
Вспомним логистическую регрессию



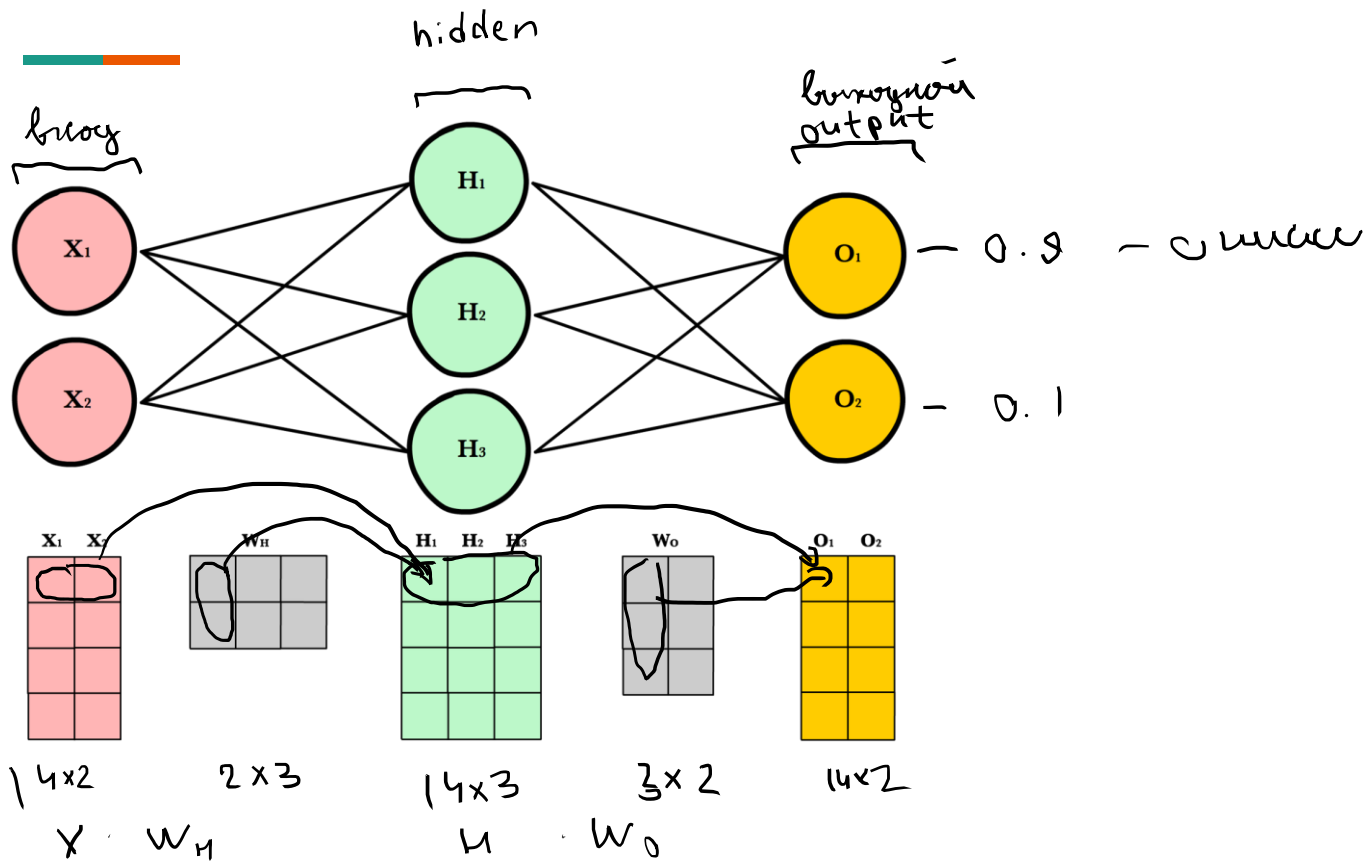
bias



А если классов больше, чем 2...



Ещё разок Summation w. w_{11}



Бинарная классификация



$$X \cdot W + b = \text{logits}$$

X

0.5	1	0	1
-----	---	---	---



$$\begin{bmatrix} 0 & -2 \\ 1 & 0 \\ 0.5 & -1 \\ -1 & 3 \end{bmatrix} + \begin{bmatrix} 0 \\ 0.5 \end{bmatrix} = \begin{bmatrix} 0 \\ 3.5 \end{bmatrix}$$

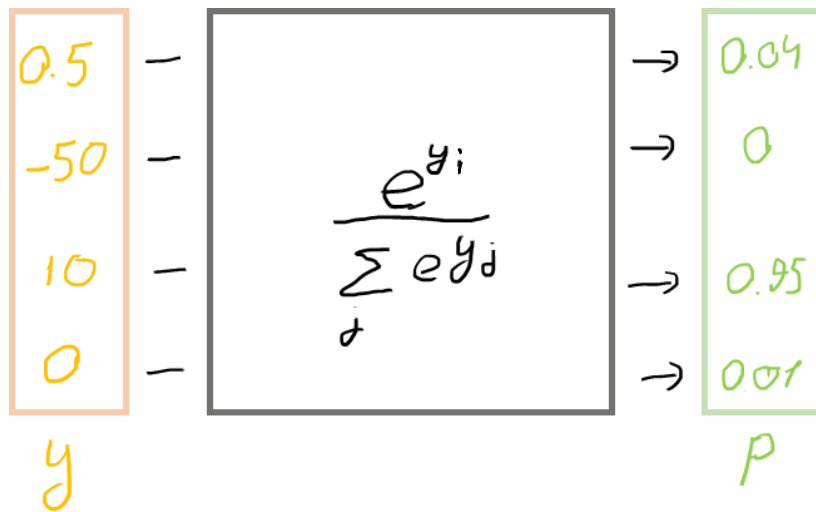
W b logits

W
4x2
число каналов
4 пружн.

Softmax



logits



Принцип максимального правдоподобия. Maximum likelihood

правдоподоб.

gt = ground truth



$$\prod_s p(c = gt_s | x_s) \rightarrow \max$$

0.9 1 1


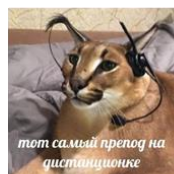
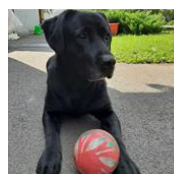
$$\frac{1}{N} \sum_s \ln p(c = gt_s | x_s) \rightarrow \max$$

↓ -1

$$- \frac{1}{N} \sum_s \ln p(c = gt_s | x_s)$$

↓
Cross Entropy



	samples	features	labels	predictions
1	 <p>100</p> <p>3</p>	<p>300000</p> <p>$x_1 x_2 \dots x_n$</p>	<p>0</p>	<p>0.9 0.09 0.01</p>
2	 <p>тот самый прелод на станции</p>	<p>...</p>	<p>1</p>	<p>0 -1 0</p>
3		<p>...</p>	<p>2</p>	<p>0 0 1</p>

А как учиться?



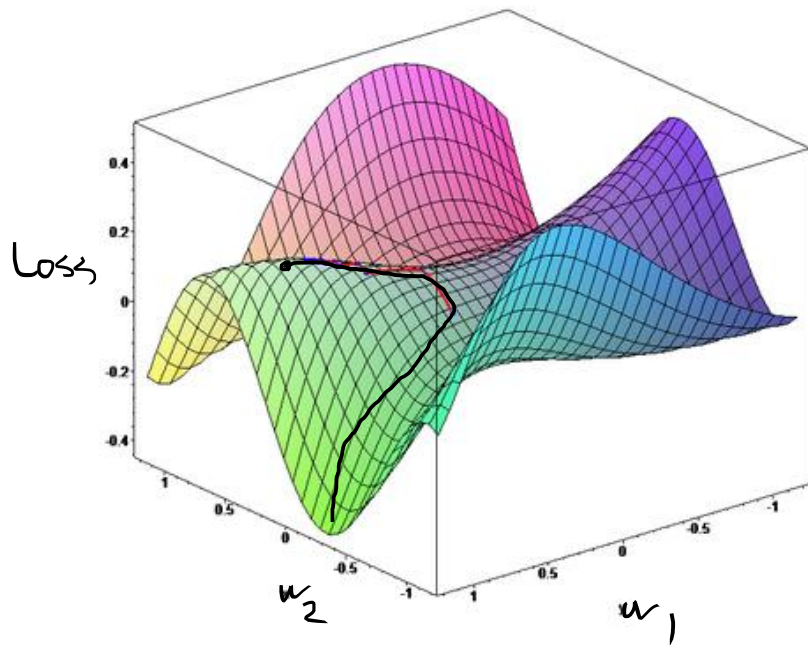
Лекции в универе похожи на просмотр Даши путешественницы: препод задаёт вопросы и несколько секунд пялится на аудиторию, а потом сам же отвечает.



А как учиться? Градиентный спуск

RAM = 8 Гб (16 Гб)

VRAM = 4 Гб



float32 - 4 байт

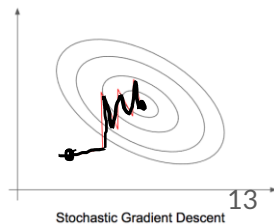
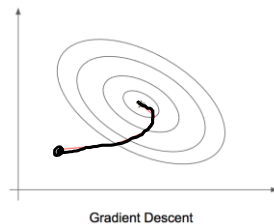
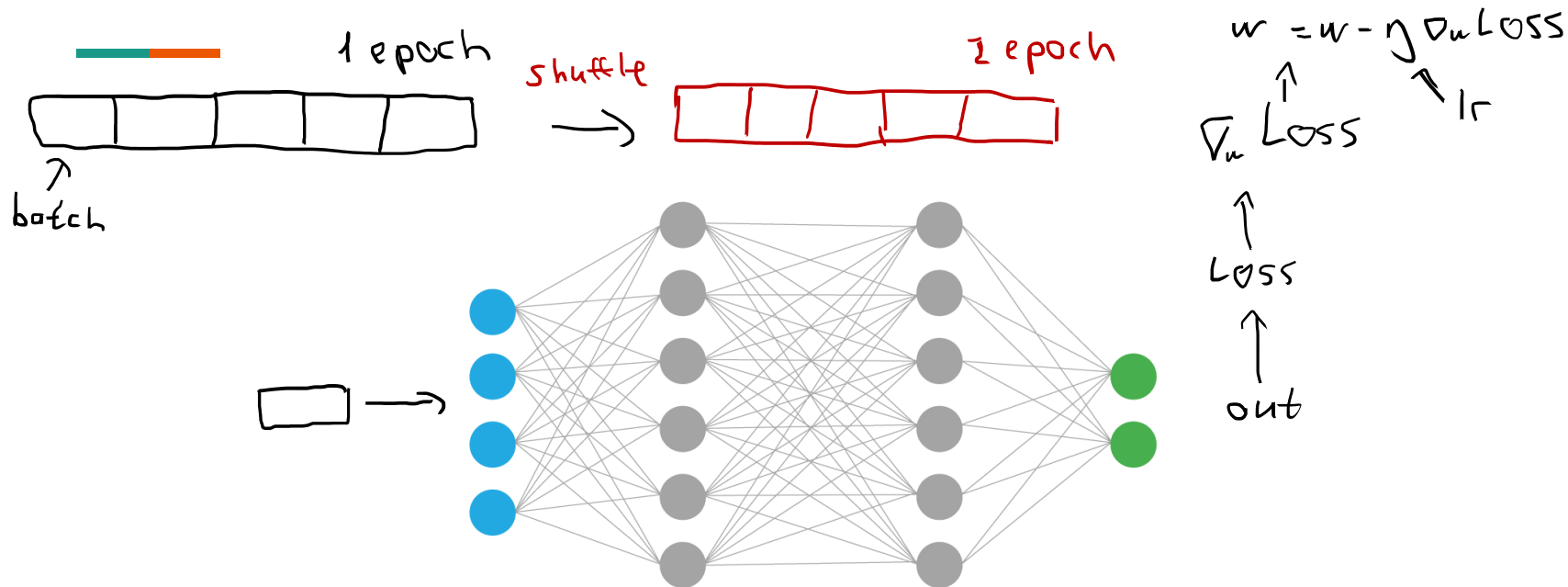
77 - 68 байт

1 000 000 - 4 мб

1 000 000 000 - 4 Гб

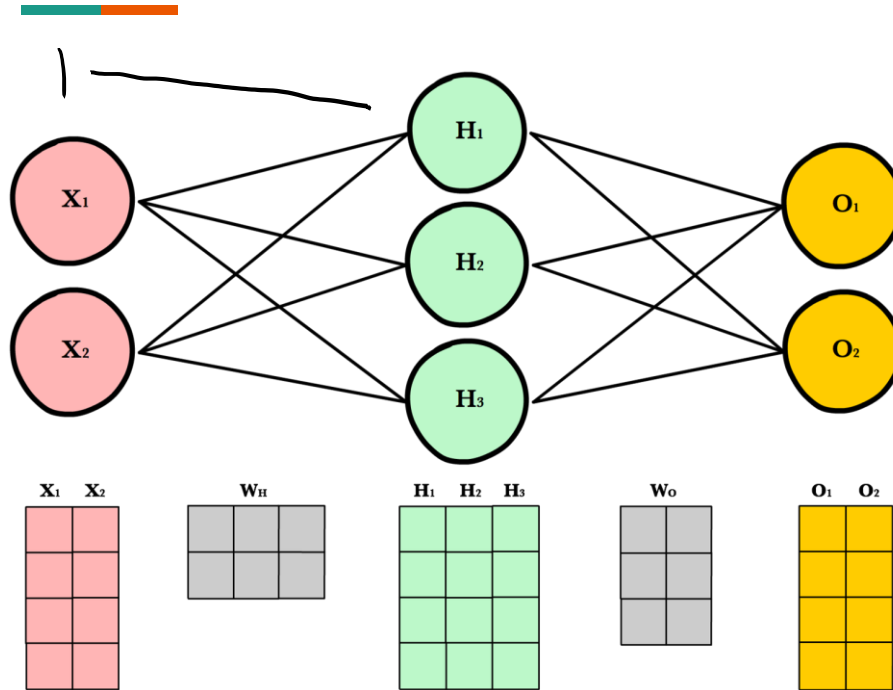
8 Гб

Стохастический градиентный спуск. Stochastic Gradient Descent (SGD)



Отдых

Будет ли работать просто так?



$$X_{4 \times 2} \cdot (W_{H_1} \cdot W_{H_2}) = O_{4 \times 2}$$

$$W_{H_1} \cdot W_{H_2} = W_{neu} \quad 2 \times 2$$

$$\Rightarrow X_{4 \times 2} \cdot W_{neu} = O_{2 \times 2}$$

Функции активации

Хотим, чтобы активация была:

Нелинейная:

Функция активации необходима для введения нелинейности в нейронные сети. Если функция активации не применяется, выходной сигнал становится простой линейной функцией. Неактивированная нейронная сеть будет действовать как линейная регрессия с ограниченной способностью к обучению:

$$\hat{y} = NN(X, W_1, \dots, W_n) = X \cdot W_1 \cdot \dots \cdot W_n = X \cdot W$$

Только нелинейные функции активации позволяют нейронным сетям решать задачи аппроксимации нелинейных функций:

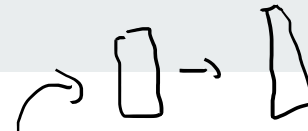
$$\hat{y} = NN(X, W_1, \dots, W_n) = \sigma(\dots \sigma(X \cdot W_1) \dots W_n) \neq X \cdot W$$

Дифференцируемая:

Функции активации должны быть дифференцируемые, то есть от них можно взять производную

Функции активации. Сигмоида

стандарти.



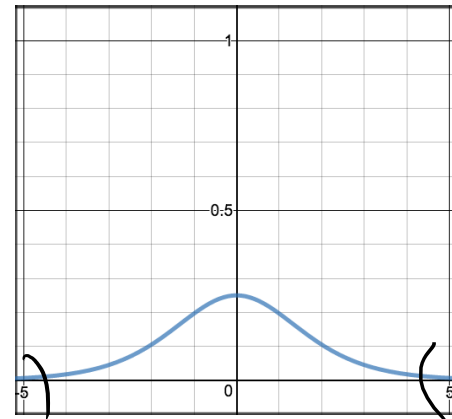
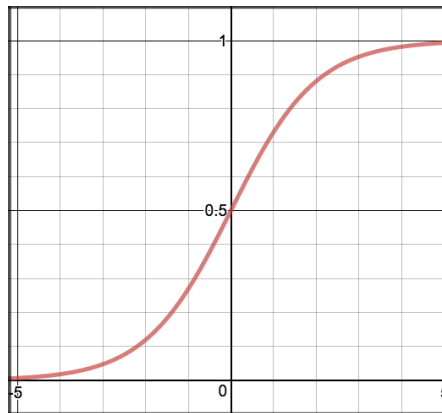
$$\frac{e^x + 1}{e^x + 1} - \frac{e^x}{e^x + 1} = \frac{1}{e^x + 1}$$

$$\sigma(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{e^x + 1}$$

$$\sigma(x)' = \frac{e^x \cdot (e^x + 1) - e^x \cdot e^x}{(e^x + 1)^2} =$$

$$= \frac{e^{2x} + e^x - e^{2x}}{(e^x + 1)^2} = \frac{e^x}{(e^x + 1)^2} =$$


$$= \frac{e^x}{e^x + 1} \cdot \frac{1}{e^x + 1} = \sigma(x) \cdot (1 - \sigma(x))$$

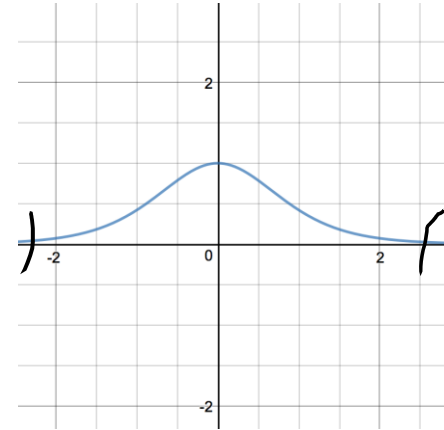
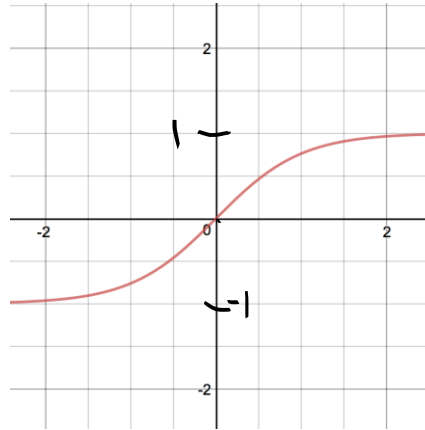


Недостатки:

1. Насыщение сигмоиды приводит к затуханию градиентов ✓
2. Выход сигмоиды не центрирован относительно нуля ✓

Функции активации. Гиперболический тангенс


$$\tanh(x) = \frac{2}{1+e^{-2x}} - 1$$



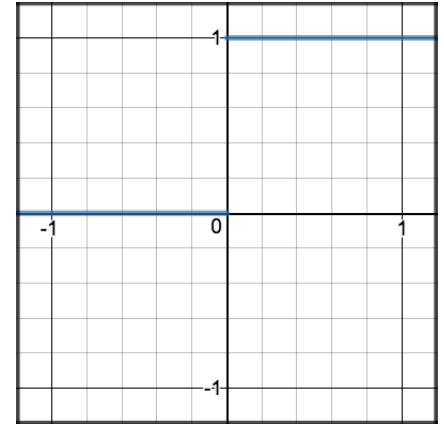
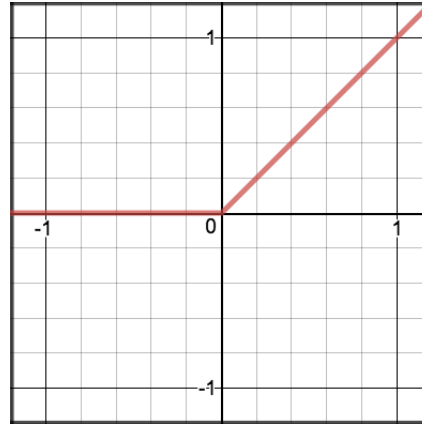
Недостатки:

1. Снова затухание градиентов

Функции активации. ReLU (rectified linear unit)



$$\text{relu}(x) = \max(0, x)$$



Недостатки:

1. Может “умереть” в области слева от 0

Функции активации. Другие функции активации

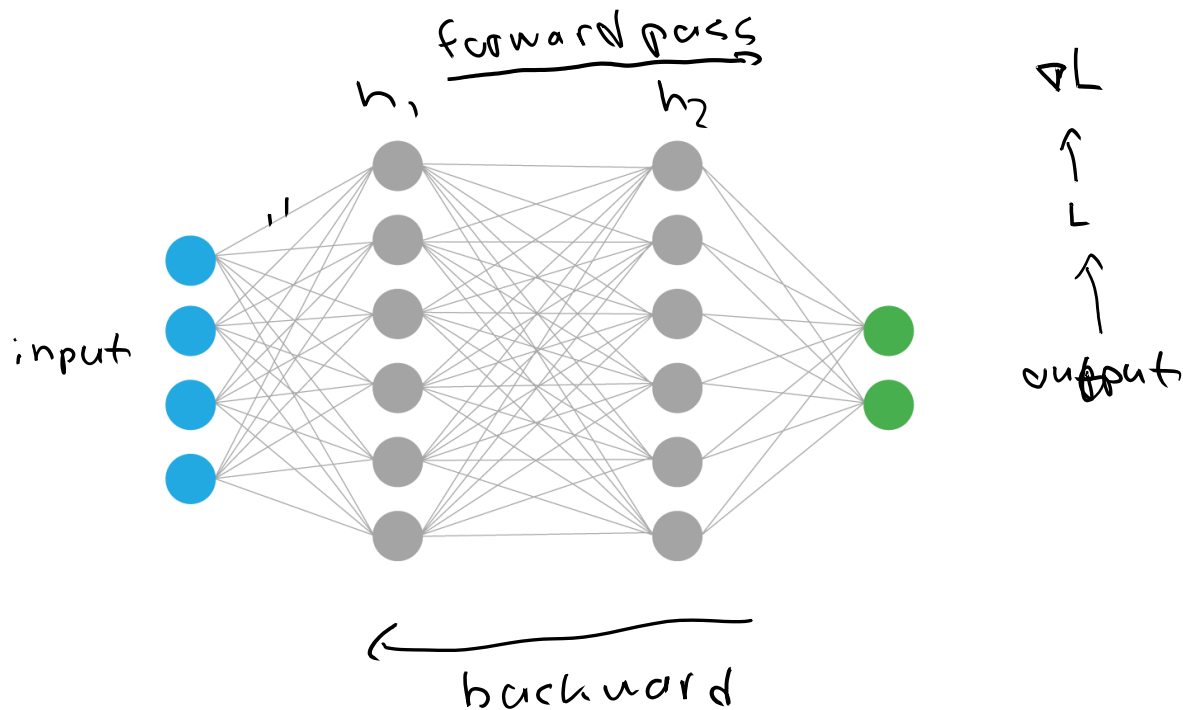


Identity	Sigmoid	TanH	ArcTan
ReLU	Leaky ReLU	Randomized ReLU	Parametric ReLU
Binary	Exponential Linear Unit	Soft Sign	Inverse Square Root Unit (ISRU)
Inverse Square Root Linear	Square Non-Linearity	Bipolar ReLU	Soft Plus

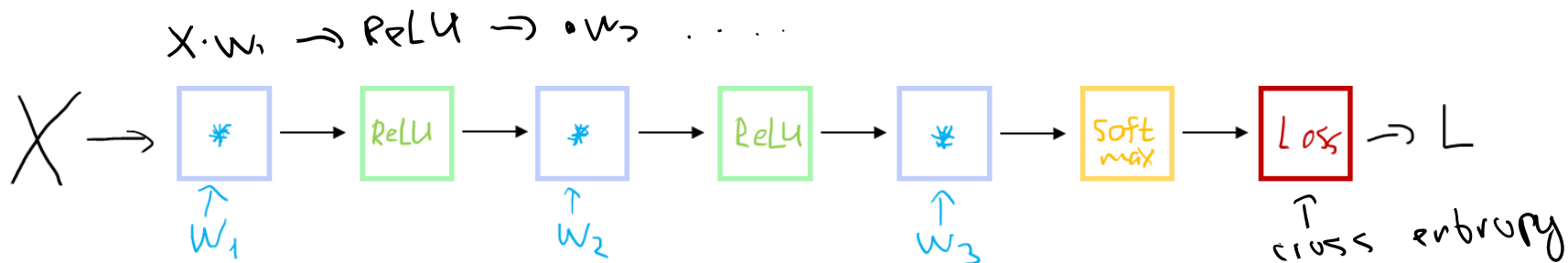
Тренировка

$$\frac{L(w + \Delta w) - L(w)}{\Delta x}$$

$$w_{i+1} = w_i - \eta \cdot \nabla_{w_i} L$$



Граф вычислений $\text{ReLU} = \max(0, x)$



Алгоритм обратного распространения ошибки. Backpropagation

Алгоритм обратного распространения ошибки позволяет находить градиенты для любого графа вычислений, если функция которую он описывает дифференцируема (каждый из узлов дифференцируемый). В его основе лежит правило взятия производной сложной функции (chain rule).

$$f(g(x))$$

$$\frac{df}{dx} = \frac{df}{dg} \cdot \frac{dg}{dx}$$

$$f(g(h(x)))$$

$$\frac{df}{dx} = \frac{df}{dg} \cdot \frac{dg}{dh} \cdot \frac{dh}{dx}$$

Тренируемся

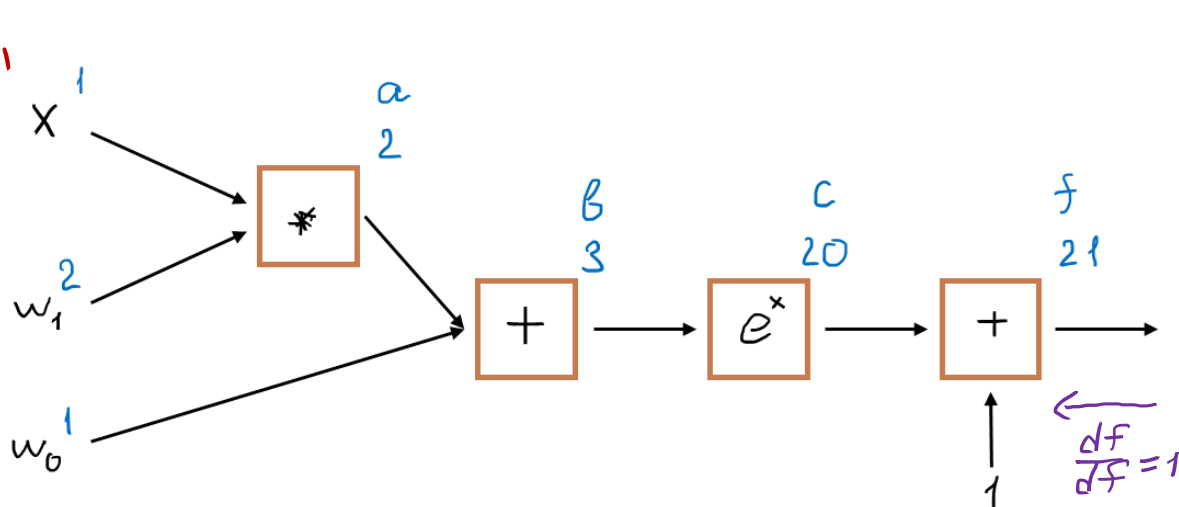
$$f = c + 1$$

$$f(g(x))$$

$$\frac{df}{dx} = \frac{df}{dg} \frac{dg}{dx}$$

$$f(x, w_0, w_1) = 1 + e^{w_0 + w_1 x}$$

$$\frac{df}{dc} = \frac{df}{dc} \cdot \frac{dc}{dc} = 1$$



Тренируемся

$$c = e^b$$

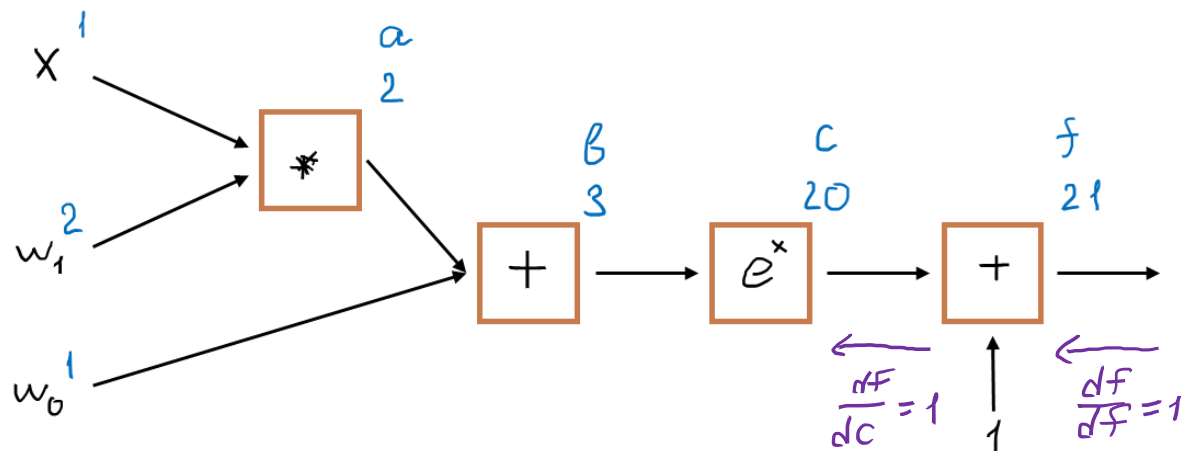
$$\frac{df}{db} = \frac{df}{dc} \cdot \frac{dc}{db}$$

$\begin{matrix} \text{"} & \text{"} \\ 1 & 20 \end{matrix}$

$$f(g(x))$$

$$\frac{df}{dx} = \frac{df}{dg} \frac{dg}{dx}$$

$$f(x, w_0, w_1) = 1 + e^{w_0 + w_1 x}$$



Тренируемся



$$f(g(x))$$

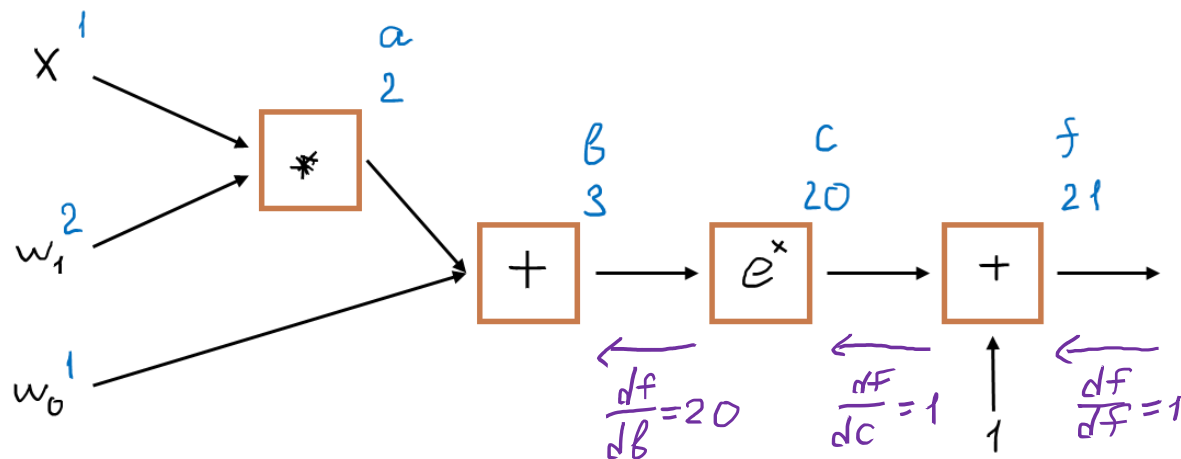
$$\frac{df}{dx} = \frac{df}{dg} \frac{dg}{dx}$$

$$f(x, w_0, w_1) = 1 + e^{w_0 + w_1 x}$$

$$b = a + w_0$$

$$\frac{df}{dw_0} = \frac{df}{db} \cdot \frac{db}{dw_0}$$

$\begin{matrix} \text{"} & \text{"} \\ 20 & 1 \end{matrix}$



Тренируемся



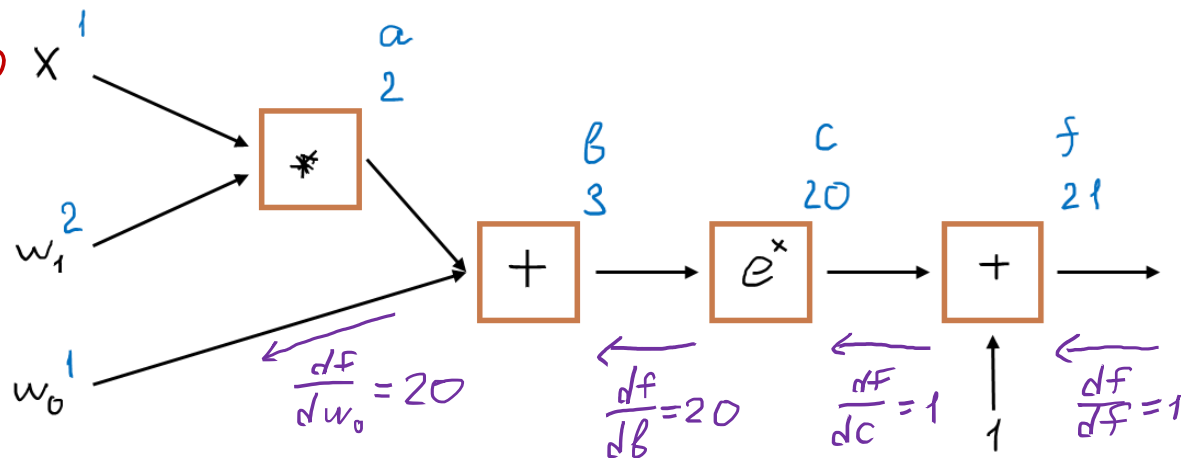
$$f(g(x))$$

$$\frac{df}{dx} = \frac{df}{dg} \frac{dg}{dx}$$

$$f(x, w_0, w_1) = 1 + e^{w_0 + w_1 x}$$

$$b = a + w_0$$

$$\frac{df}{da} = \frac{df}{db} \cdot \frac{db}{da} = 20 \times 1 = 20$$



Тренируемся

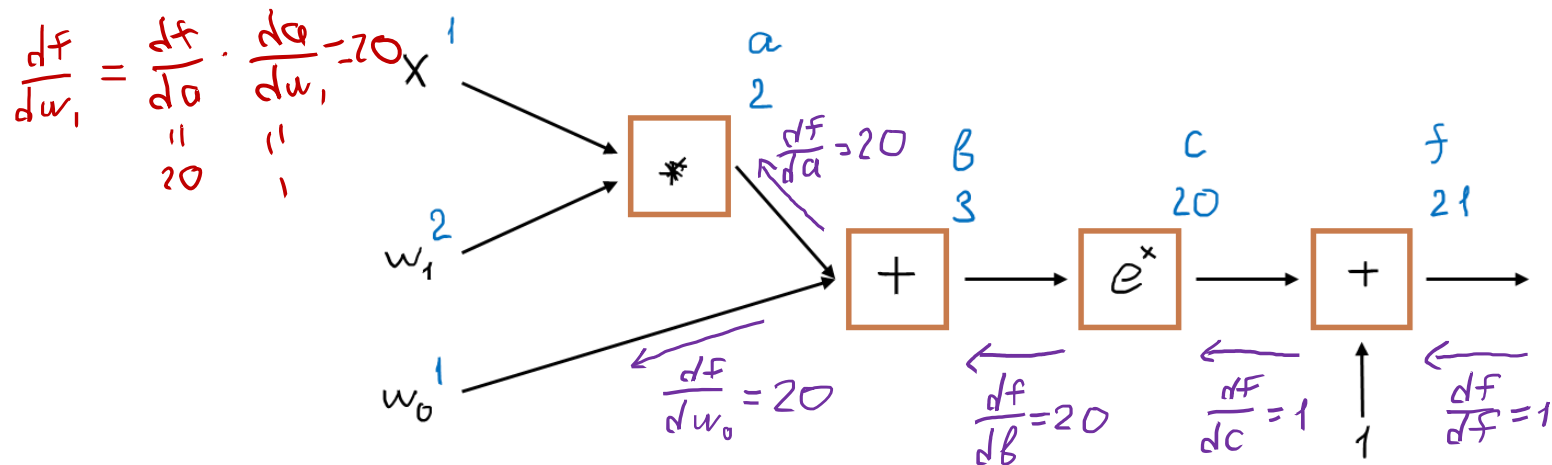


$$f(g(x))$$

$$\frac{df}{dx} = \frac{df}{dg} \frac{dg}{dx}$$

$$f(x, w_0, w_1) = 1 + e^{w_0 + w_1 x}$$

$$a = w_1 \cdot x$$



Тренируемся

$$a = w_1 \cdot x$$

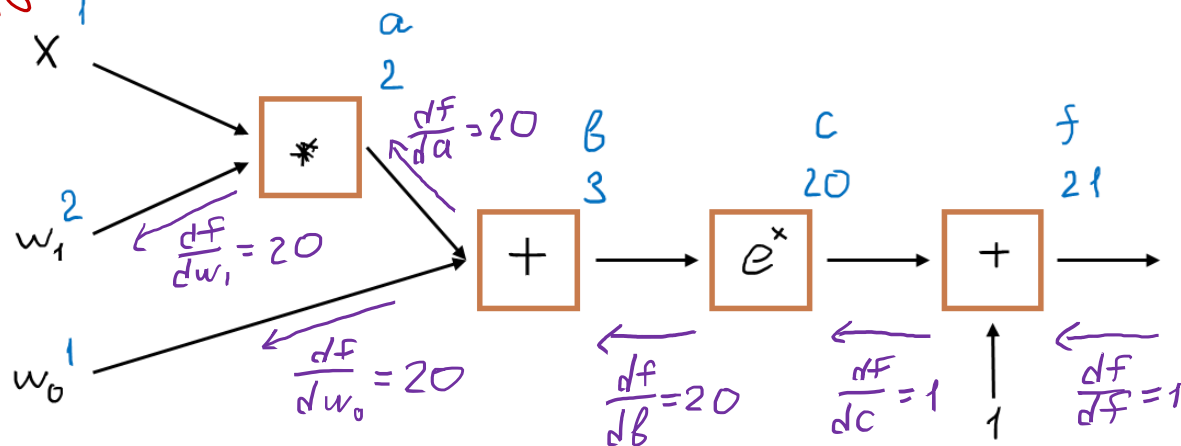
$$\frac{df}{dx} = \frac{df}{da} \cdot \frac{da}{dx} = 40$$

Handwritten calculation: $\frac{df}{dx} = \frac{df}{da} \cdot \frac{da}{dx} = 40$. The values 20 and 2 are written below the fractions, indicating $\frac{df}{da} = 20$ and $\frac{da}{dx} = 2$.

$$f(g(x))$$

$$\frac{df}{dx} = \frac{df}{dg} \frac{dg}{dx}$$

$$f(x, w_0, w_1) = 1 + e^{w_0 + w_1 x}$$

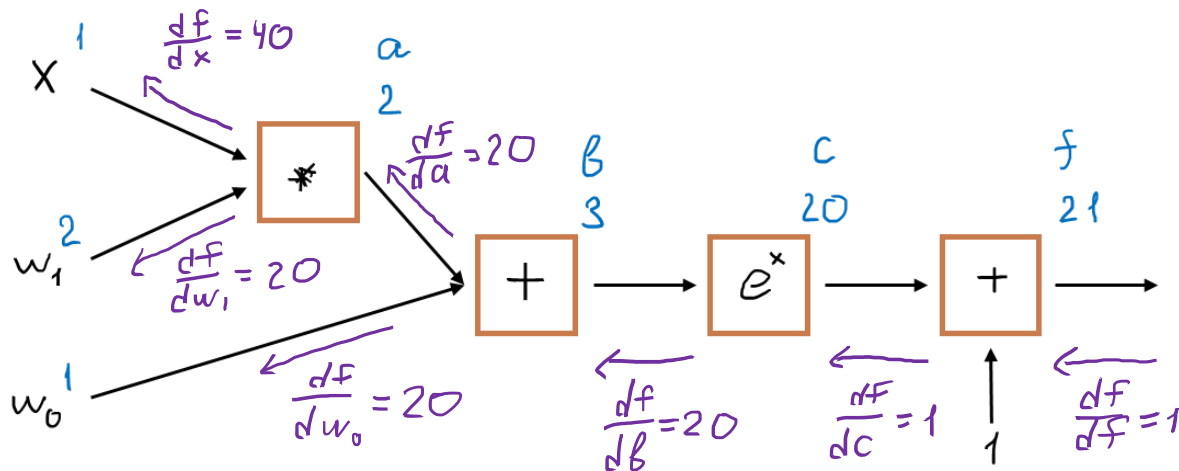


Тренируемся

$$f(g(x))$$

$$\frac{df}{dx} = \frac{df}{dg} \frac{dg}{dx}$$

$$f(x, w_0, w_1) = 1 + e^{w_0 + w_1 x}$$



Тренируемся



$$f(g(x))$$

$$\frac{df}{dx} = \frac{df}{dg} \frac{dg}{dx}$$

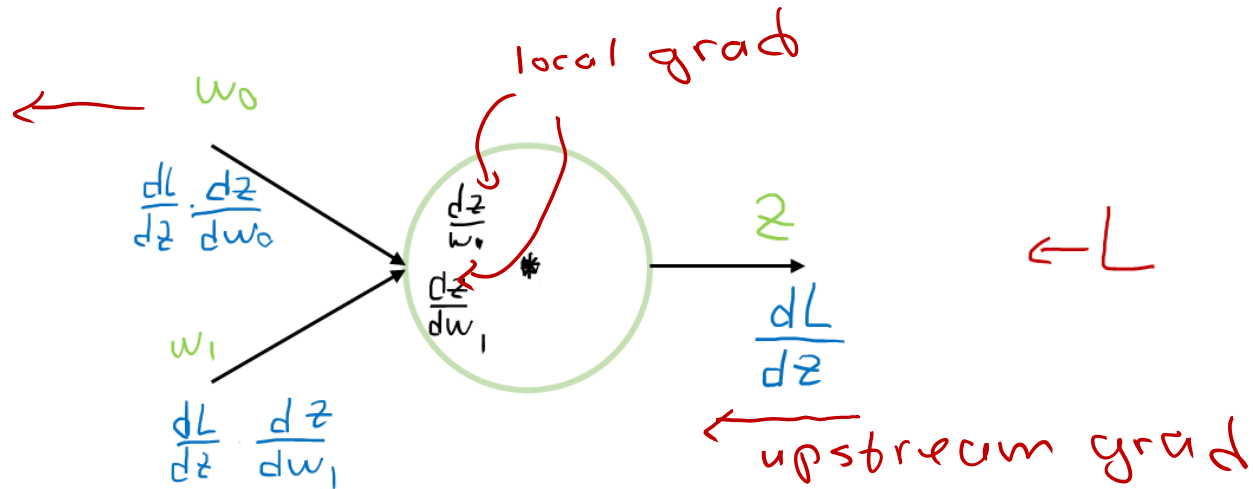
$$f(x, w_0, w_1) = 1 + e^{w_0 + w_1 x}$$

А зачем мы все это считали?

$$\frac{df}{dx}, \frac{df}{dw_0}, \frac{df}{dw_1},$$

$$x = x - \eta \cdot \frac{df}{dx}$$

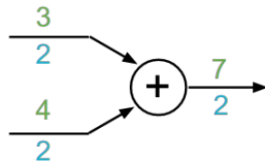
Общая схема вычисления градиента



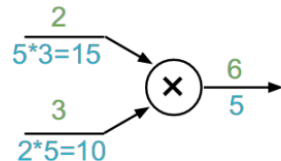
Уточнения



add gate: gradient distributor



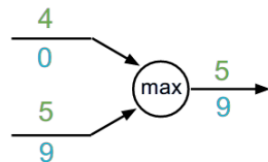
mul gate: "swap multiplier"



copy gate: gradient adder



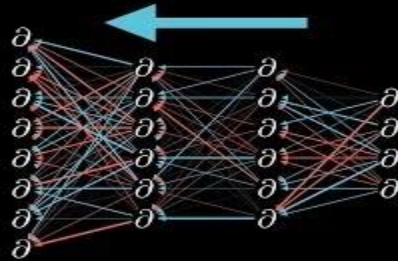
max gate: gradient router



Если ничего не понятно...



Backpropagation calculus



знаешь почему ты
так сильно устаёшь к концу дня?
потому что ты весь день был
замечательным котёночком,
а это тяжелый труд

