

# Оптимизация в нейронных сетях.

## База

# Обучение NN и параллельные вычисления

$$L(W, X, y) = \frac{1}{N} \sum_{i=1}^N l(W, x_i, y_i) \rightarrow \min_{W \in \mathbb{R}^p}$$

# Обучение NN и параллельные вычисления

Функция потерь  
(меньше - лучше)

$$L(W, X, y) = \frac{1}{N} \sum_{i=1}^N l(W, x_i, y_i) \rightarrow \min_{W \in \mathbb{R}^p}$$

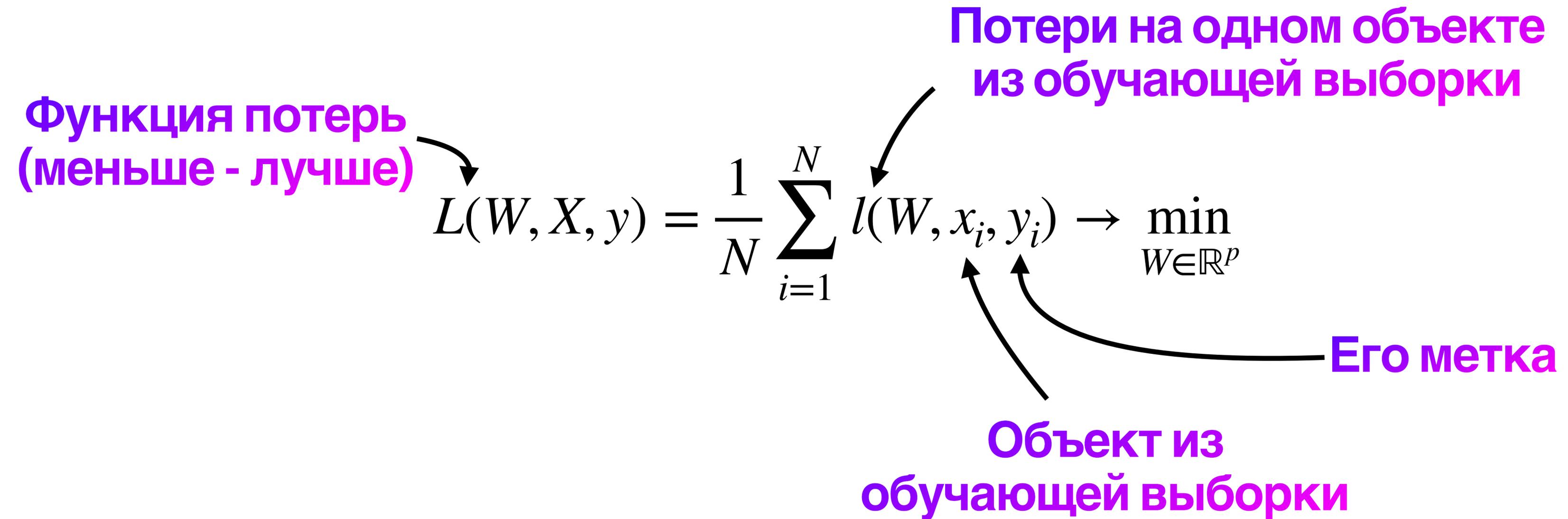
# Обучение NN и параллельные вычисления

Функция потерь  
(меньше - лучше)

$$L(W, X, y) = \frac{1}{N} \sum_{i=1}^N l(W, x_i, y_i) \rightarrow \min_{W \in \mathbb{R}^p}$$

Потери на одном объекте  
из обучающей выборки

# Обучение NN и параллельные вычисления



# Обучение NN и параллельные вычисления

Размер обучающей выборки

ImageNet  $\approx 1.4 \cdot 10^7$

WikiText  $\approx 10^8$

Потери на одном объекте  
из обучающей выборки

Функция потерь  
(меньше - лучше)

$$L(W, X, y) = \frac{1}{N} \sum_{i=1}^N l(W, x_i, y_i) \rightarrow \min_{W \in \mathbb{R}^p}$$

Его метка

Объект из  
обучающей выборки

# Обучение NN и параллельные вычисления

Размер обучающей выборки

ImageNet  $\approx 1.4 \cdot 10^7$

WikiText  $\approx 10^8$

Потери на одном объекте из обучающей выборки

Функция потерь (меньше - лучше)

$$L(W, X, y) = \frac{1}{N} \sum_{i=1}^N l(W, x_i, y_i) \rightarrow \min_{W \in \mathbb{R}^p}$$

Веса модели, которые нужно подобрать  
НО КАК?

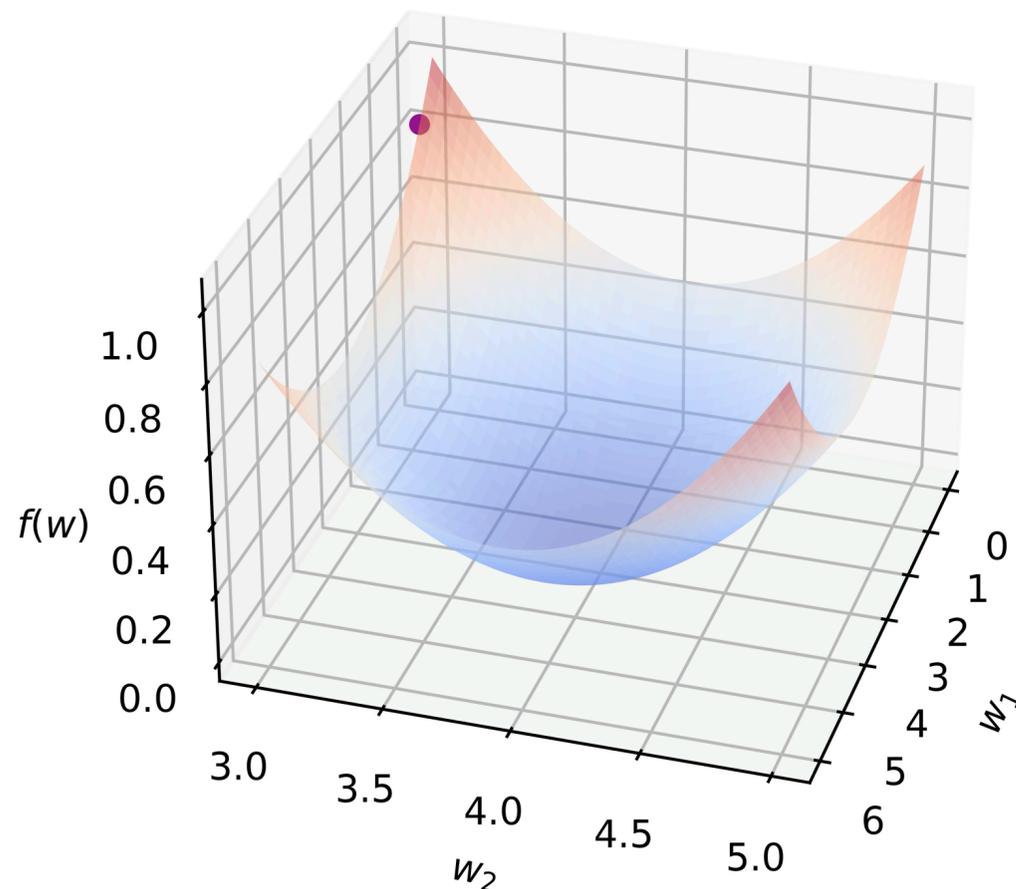
Его метка  
Объект из обучающей выборки

# Обучение NN и параллельные вычисления

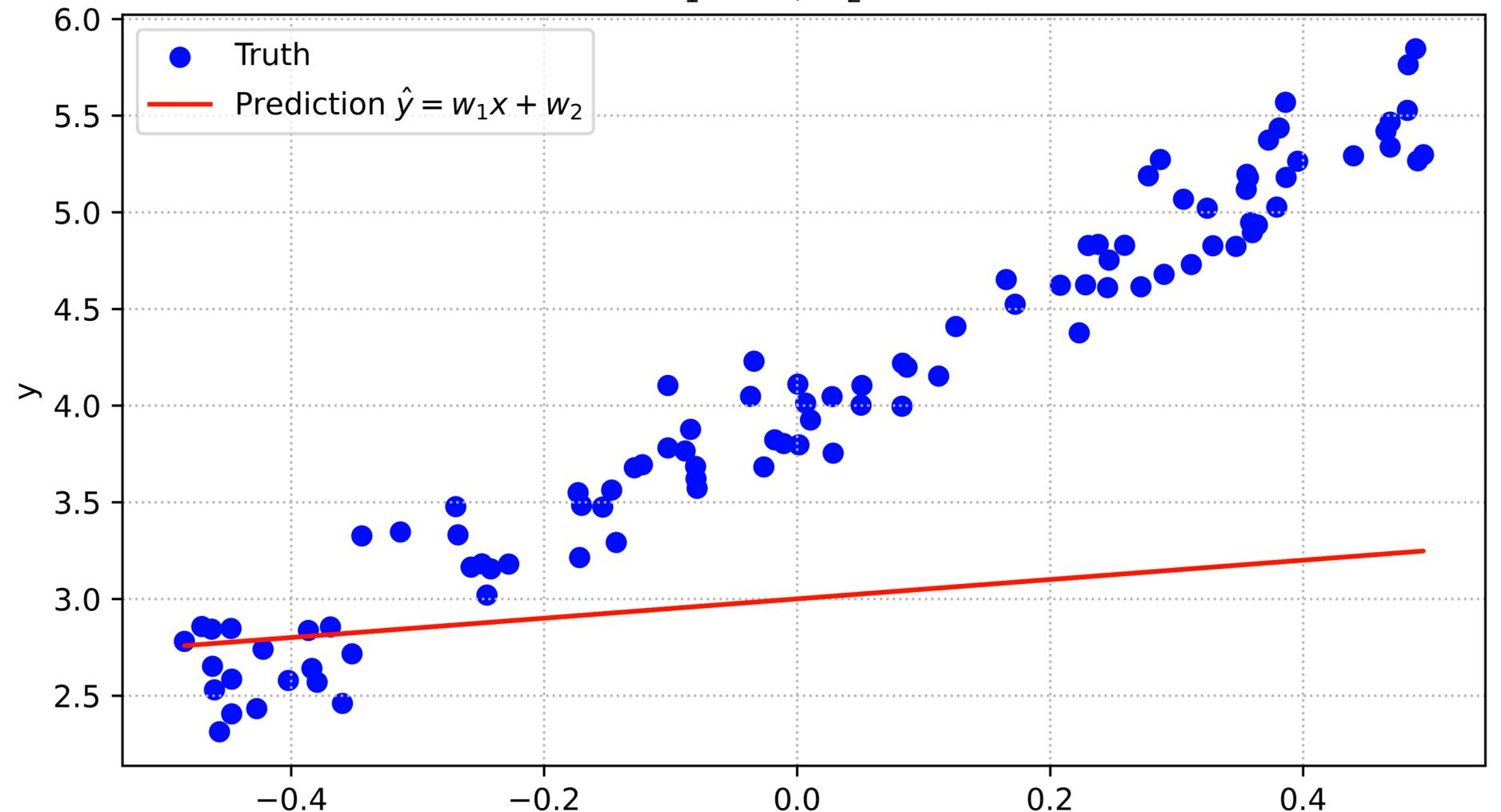
## Метод градиентного спуска (GD)

$$W_{k+1} = W_k - \alpha \nabla_W L(W_k)$$

Loss value 0.87



$w_1$  0.50,  $w_2$  3.00



# Обучение NN и параллельные вычисления

Метод градиентного спуска (GD)

$$W_{k+1} = W_k - \alpha \nabla_W L(W_k)$$

$$L(W, X, y) = \frac{1}{N} \sum_{i=1}^N l(W, x_i, y_i) \rightarrow \min_{W \in \mathbb{R}^p}$$

# Обучение NN и параллельные вычисления

Метод градиентного спуска (GD)

$$W_{k+1} = W_k - \alpha \nabla_W L(W_k)$$

$$L(W, X, y) = \frac{1}{N} \sum_{i=1}^N l(W, x_i, y_i) \rightarrow \min_{W \in \mathbb{R}^p}$$

$$\nabla_W L(W_k) = \frac{1}{N} \sum_{i=1}^N \nabla_W l(W_k, x_i, y_i)$$

# Обучение NN и параллельные вычисления

Метод градиентного спуска (GD)

$$W_{k+1} = W_k - \alpha \nabla_W L(W_k)$$

$$L(W, X, y) = \frac{1}{N} \sum_{i=1}^N l(W, x_i, y_i) \rightarrow \min_{W \in \mathbb{R}^p}$$

$$\nabla_W L(W_k) = \frac{1}{N} \sum_{i=1}^N \nabla_W l(W_k, x_i, y_i)$$

Тяжело считать при большом N 🤔

# Обучение NN и параллельные вычисления

Метод стохастического  
градиентного спуска (SGD)

$$W_{k+1} = W_k - \alpha g_k$$

$$\nabla_W L(W_k) = \frac{1}{N} \sum_{i=1}^N \nabla_W l(W_k, x_i, y_i)$$

# Обучение NN и параллельные вычисления

Метод стохастического  
градиентного спуска (SGD)

$$W_{k+1} = W_k - \alpha g_k$$

$$\nabla_W L(W_k) = \frac{1}{N} \sum_{i=1}^N \nabla_W l(W_k, x_i, y_i)$$

$$g_k = \frac{1}{b} \sum_{i=1}^b \nabla_W l(W_k, x_{j_i}, y_{j_i})$$

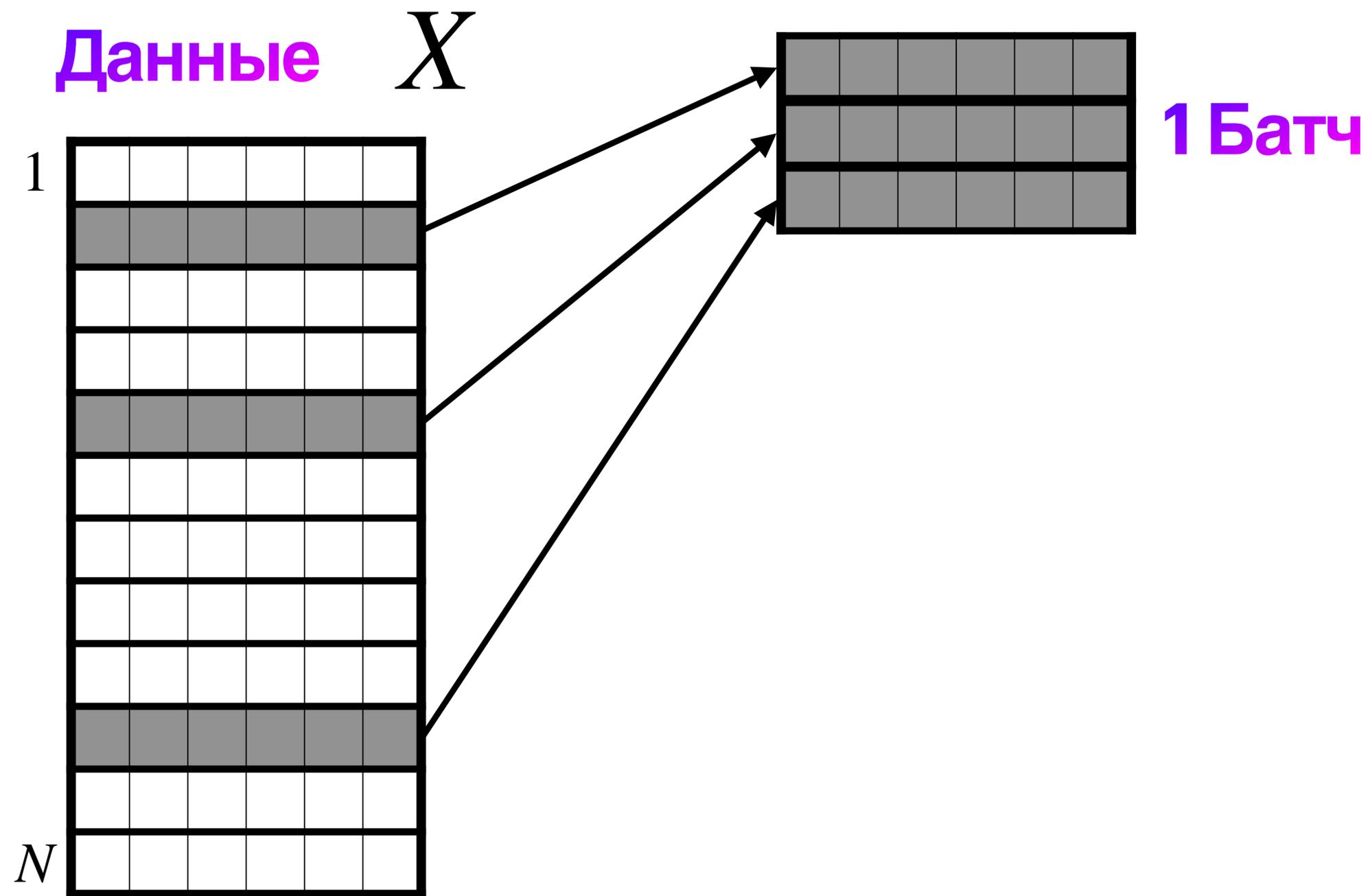
$$b \ll N$$

# Обучение NN и параллельные вычисления

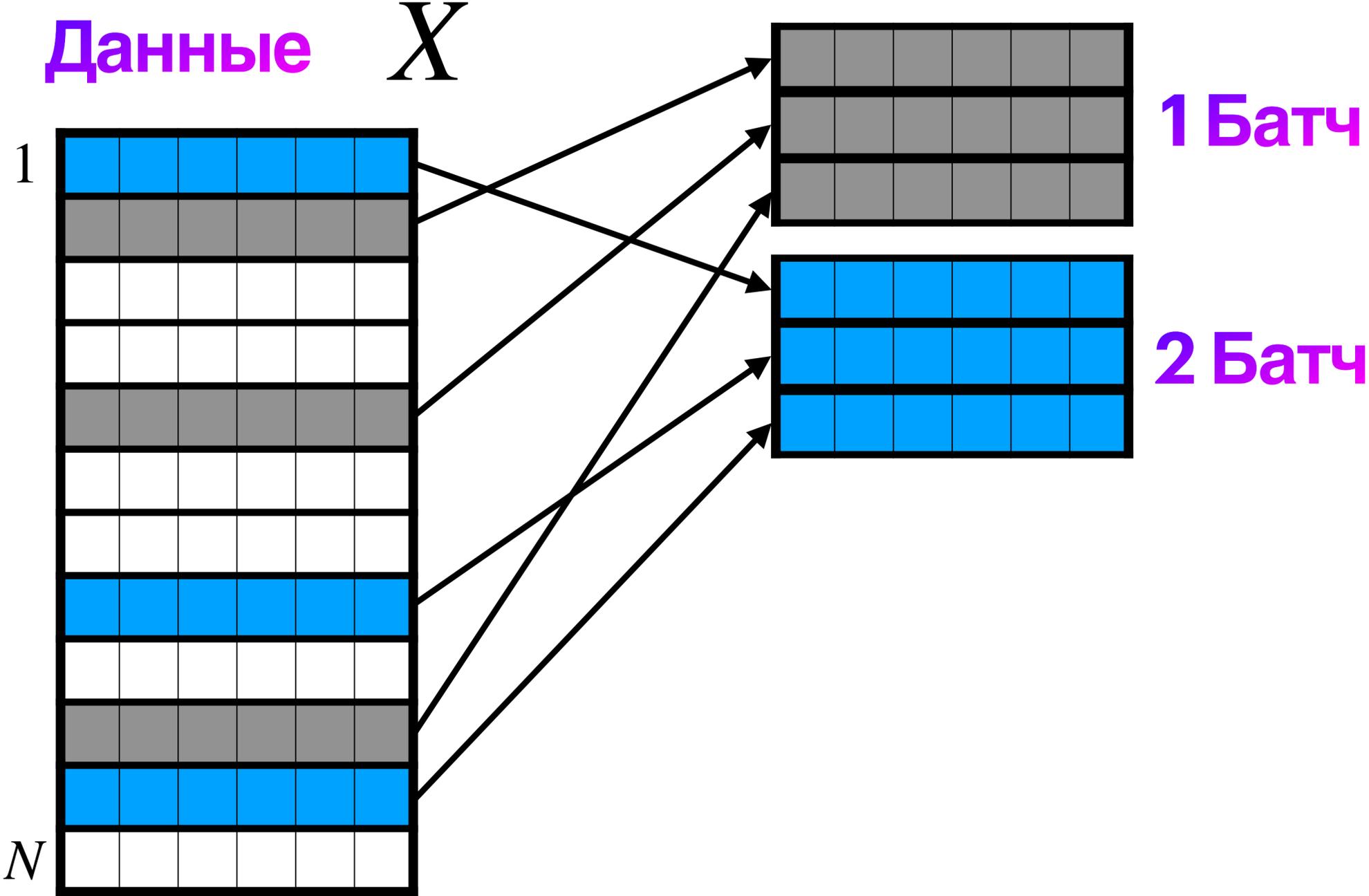
Данные  $X$

1						
$N$						

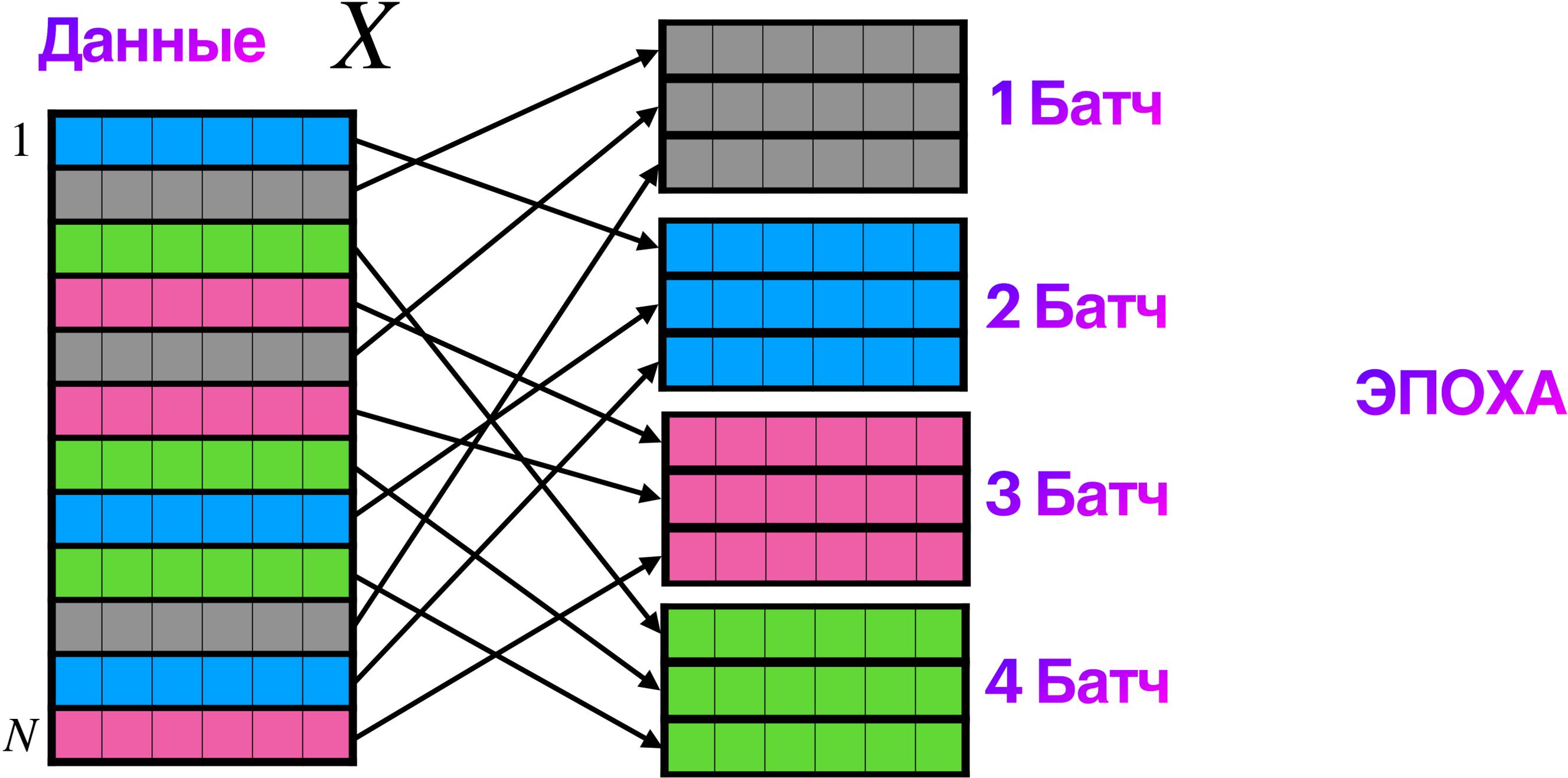
# Обучение NN и параллельные вычисления



# Обучение NN и параллельные вычисления



# Обучение NN и параллельные вычисления



# Свое железо для обучения - дорого

	ARCH	VRAM	FP32 TFLOPS	ENERGY	ENPRICE	PPRICE
RTX 2080 Ti	Turing	11	13.5	250	500 €	450 €
RTX 2080 Super	Turing	8	11.2	215	430 €	350 €
RTX 2070 Super	Turing	8	9.1	175	350 €	300 €
RTX 4090	Ada Lovelace	24	82.6	450	900 €	1,900 €
RTX 4070 Ti	Ada Lovelace	12	40.1	285	570 €	1,000 €
RTX 3070 Ti	Ampere	8	21.75	290	580 €	600 €
RTX 4080	Ada Lovelace	16	48.7	320	640 €	1,300 €
RTX 3080 Ti	Ampere	12	34.1	350	700 €	1,100 €
RTX 3090 Ti	Ampere	24	40	450	900 €	1,500 €
RTX A4000	Ampere	16	19.1	140	280 €	1,000 €
RTX 4000 Ada	Ada Lovelace	20	19.2	70	140 €	1,250 €
RTX A5000	Ampere	24	27.8	230	460 €	2,500 €
RTX 6000 Ada	Ada Lovelace	48	91.1	300	600 €	7,000 €
L40	Ada Lovelace	48	90.5	300	600 €	9,000 €
RTX A6000	Ampere	48	38.7	300	600 €	6,000 €
A100	Ampere	40	19.5	250	500 €	11,000 €
H100	Hopper	80	24.1	350	700 €	35,000 €

ENPRICE - оценка стоимости работы GPU на 100% в течение года



Сайт

Таблица ценами на GPU и сравнением

**Теория**



# Стох. градиент - несмещенная оценка настоящего градиента

$$\mathbb{E}g_k = \nabla_W L(W_k)$$

# SGD с постоянным шагом не сходится для выпуклой задачи

“Сходимость” для сильно выпуклой задачи:

$$\mathbb{E}[\|x_k - x^*\|^2] \leq (1 - 2\alpha_k\mu)^k R^2 + \frac{\alpha B^2}{2\mu},$$

# SGD в отличие от GD может “выпрыгивать” из локальных минимумов

# Как работает 🧠 ?

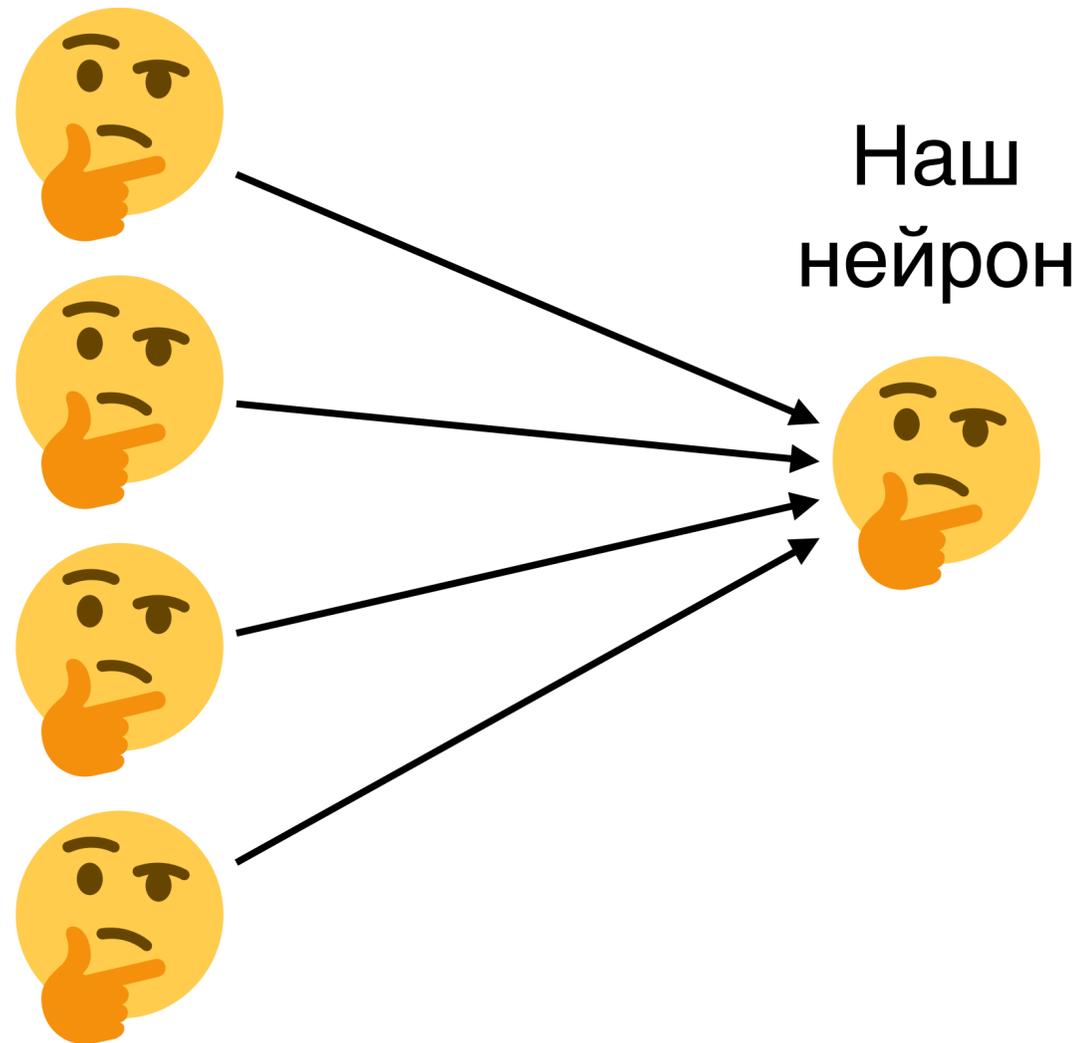


$8.6 \times 10^{10}$  нейронов в мозге человека

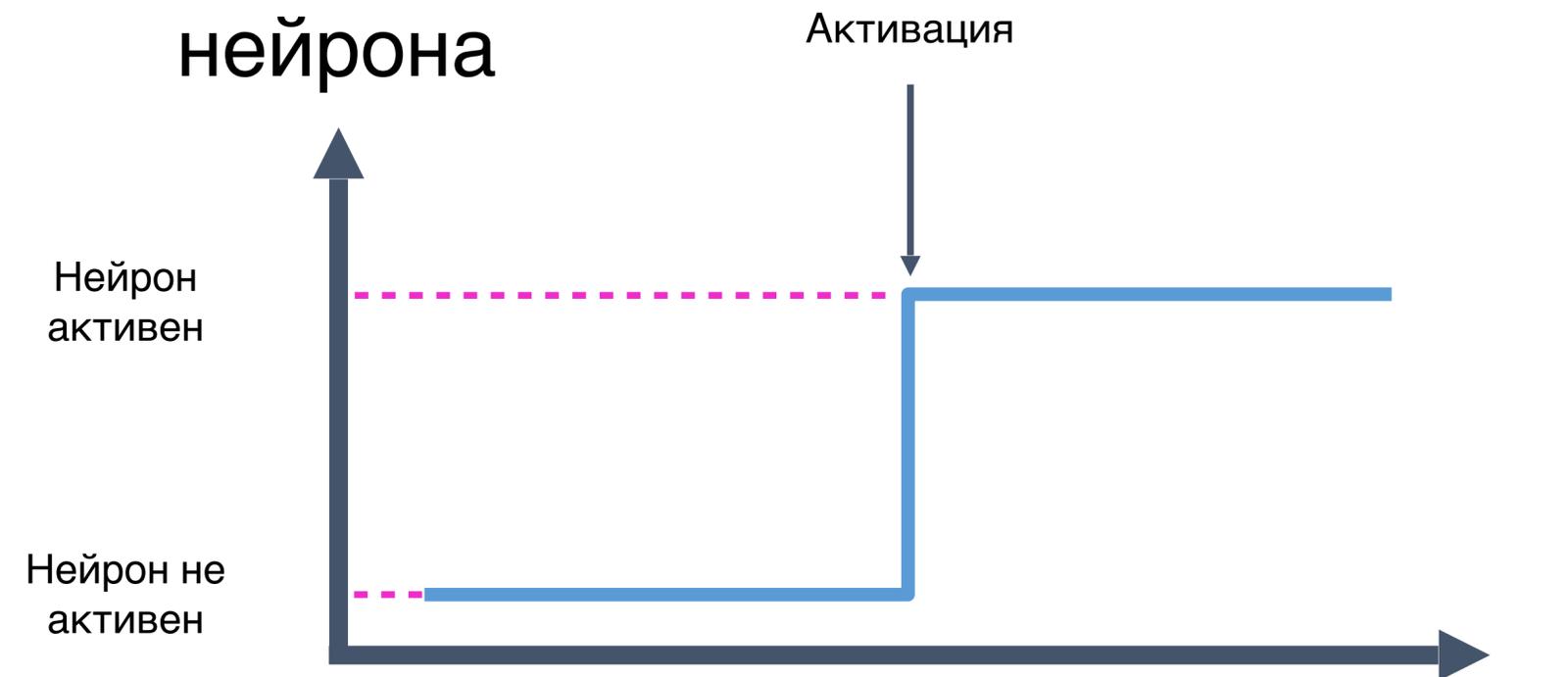
$\sim 1.5 \times 10^{14}$  связей между ними

# Как работает 🧠 ?

Другие  
нейроны



Реакция  
нейрона



Сумма входных сигналов  
от других нейронов

# Как работает 🤖🧠 ?

Другие  
нейроны



$x_1$



$x_2$



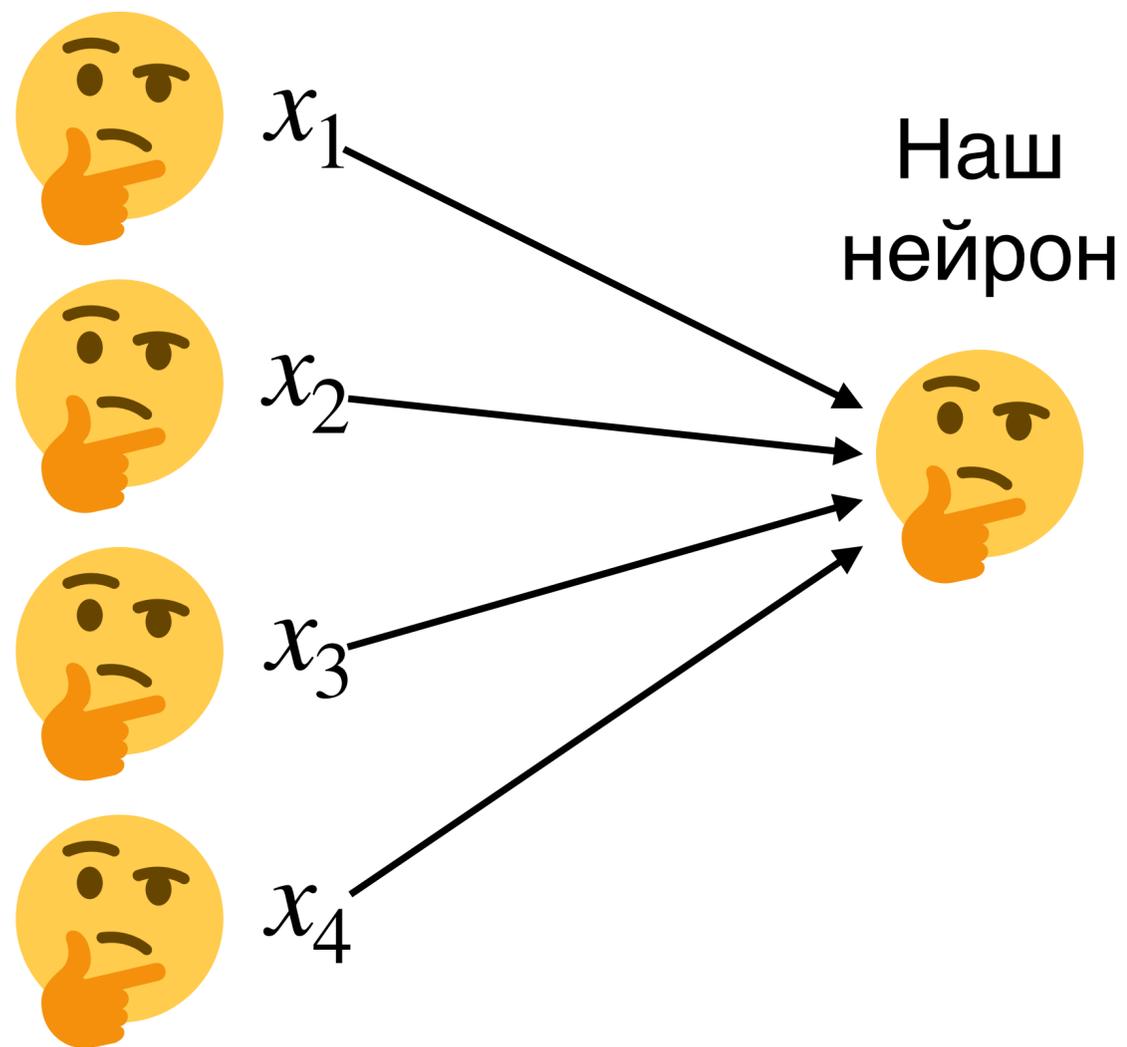
$x_3$



$x_4$

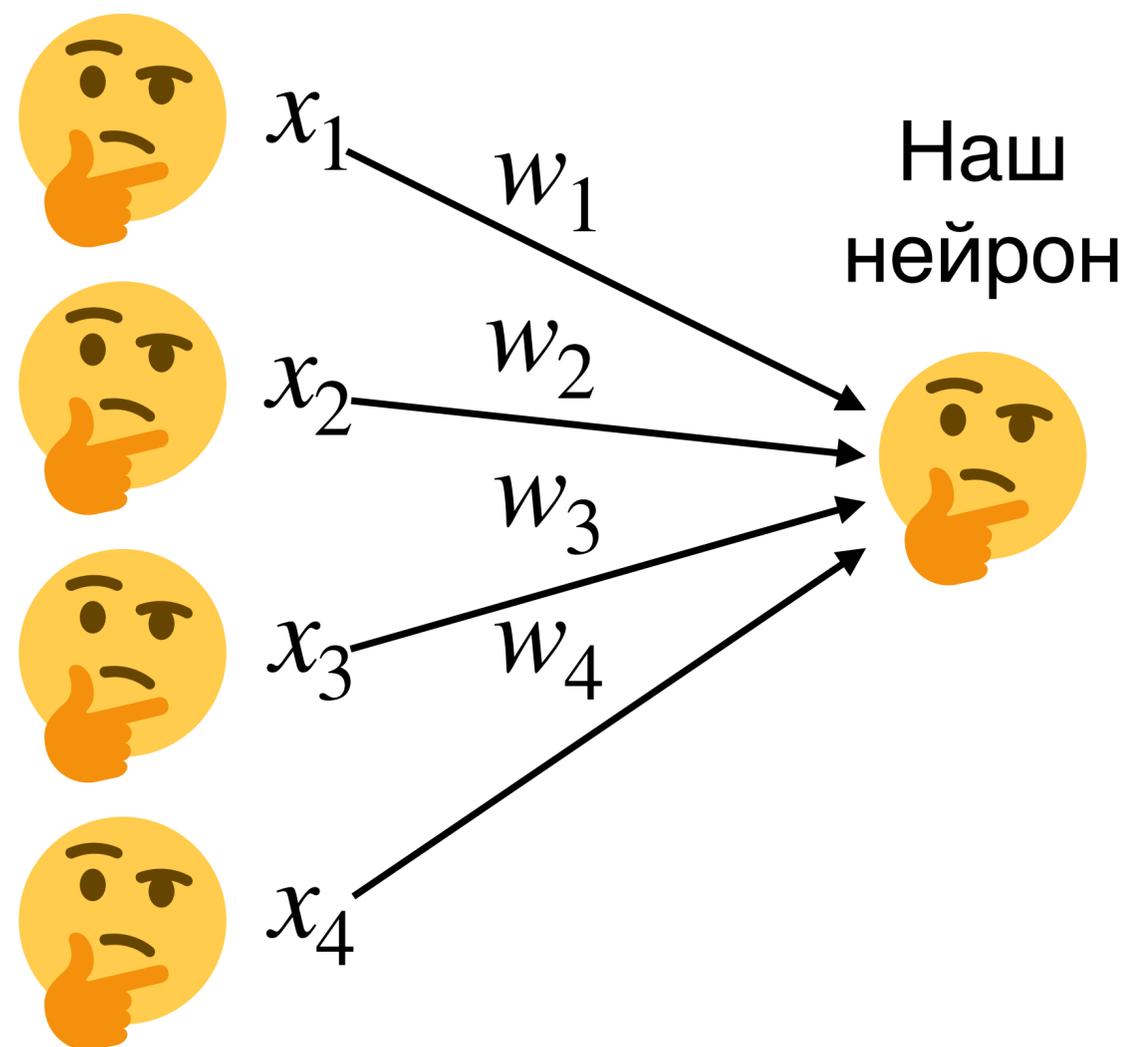
# Как работает 🤖🧠 ?

Другие  
нейроны



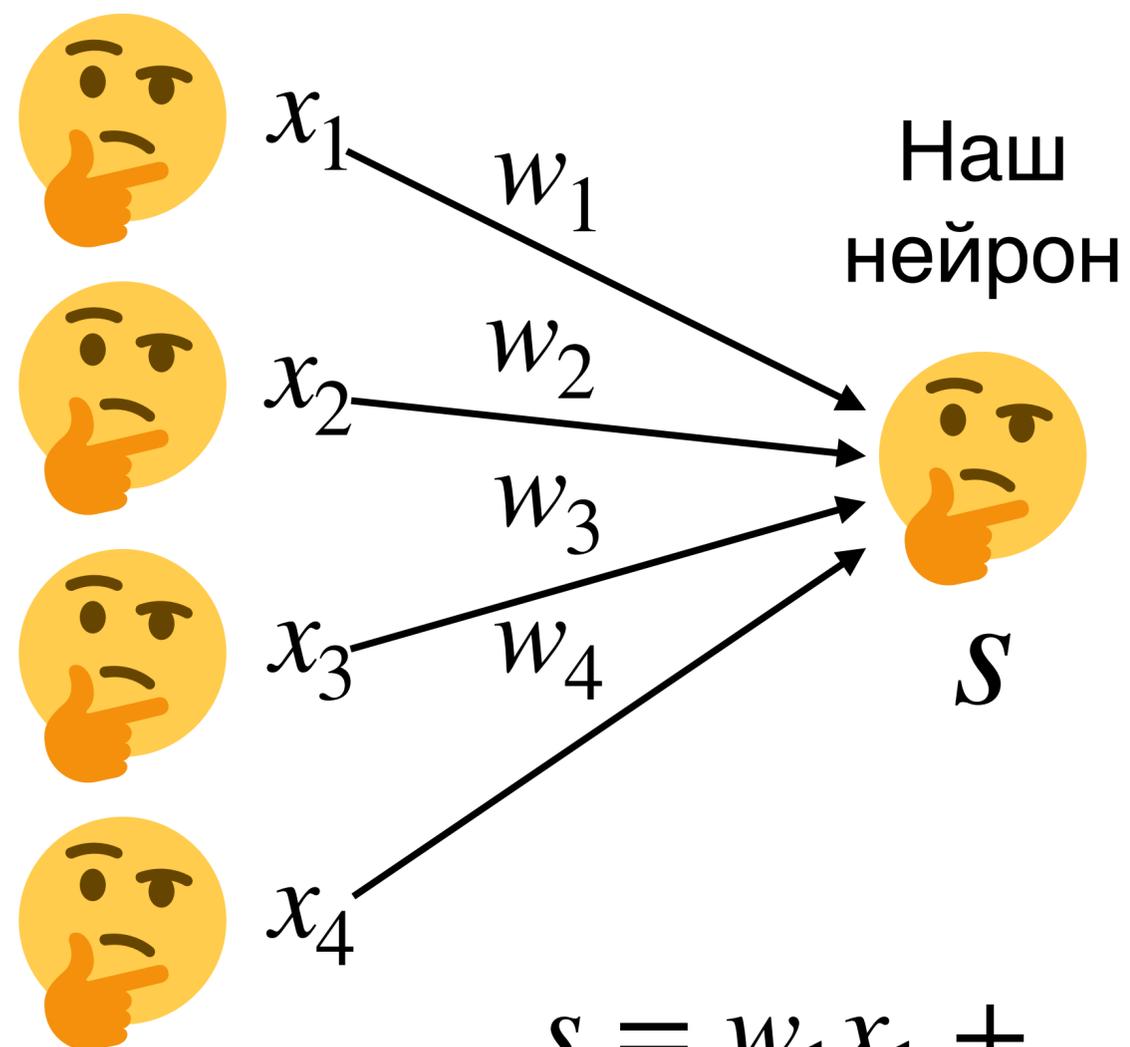
# Как работает 🤖🧠 ?

Другие  
нейроны



# Как работает 🤖🧠 ?

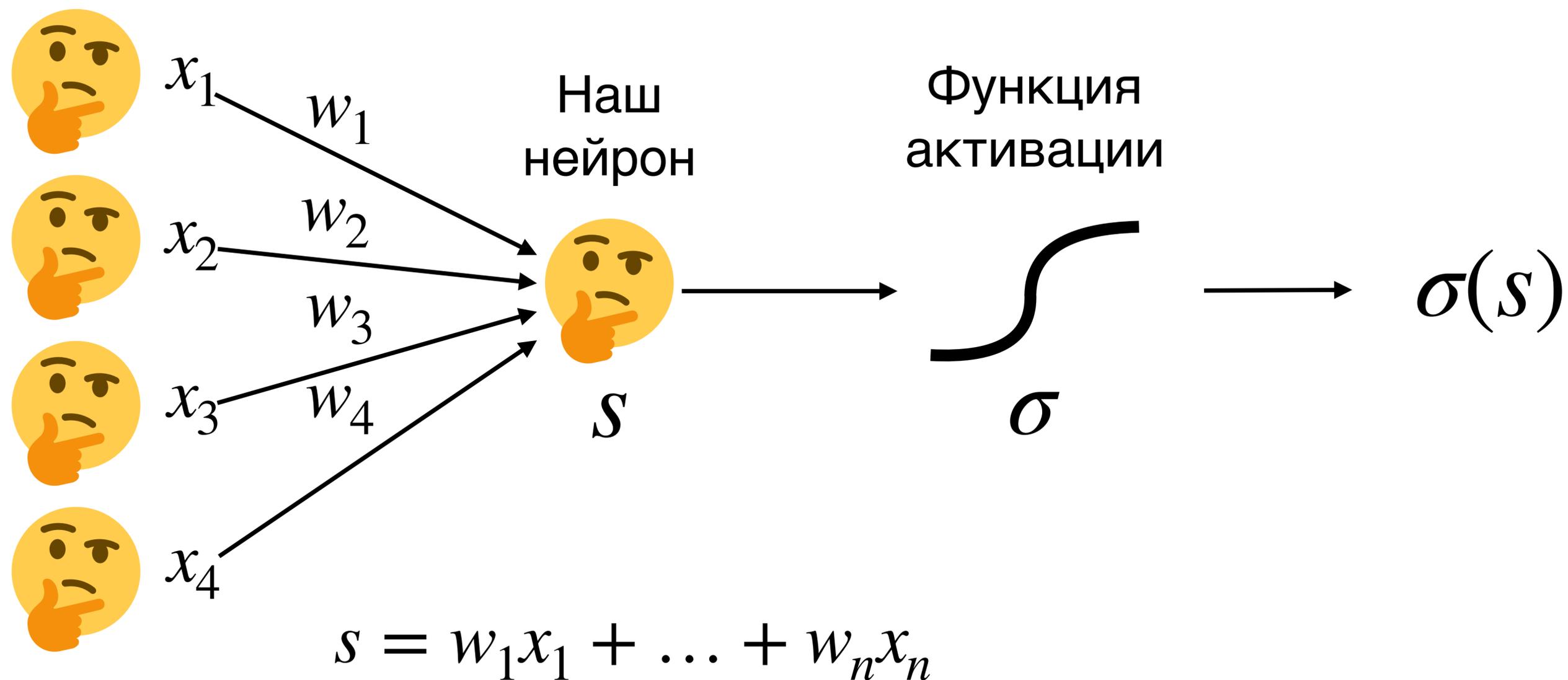
Другие  
нейроны



$$S = w_1x_1 + \dots + w_nx_n$$

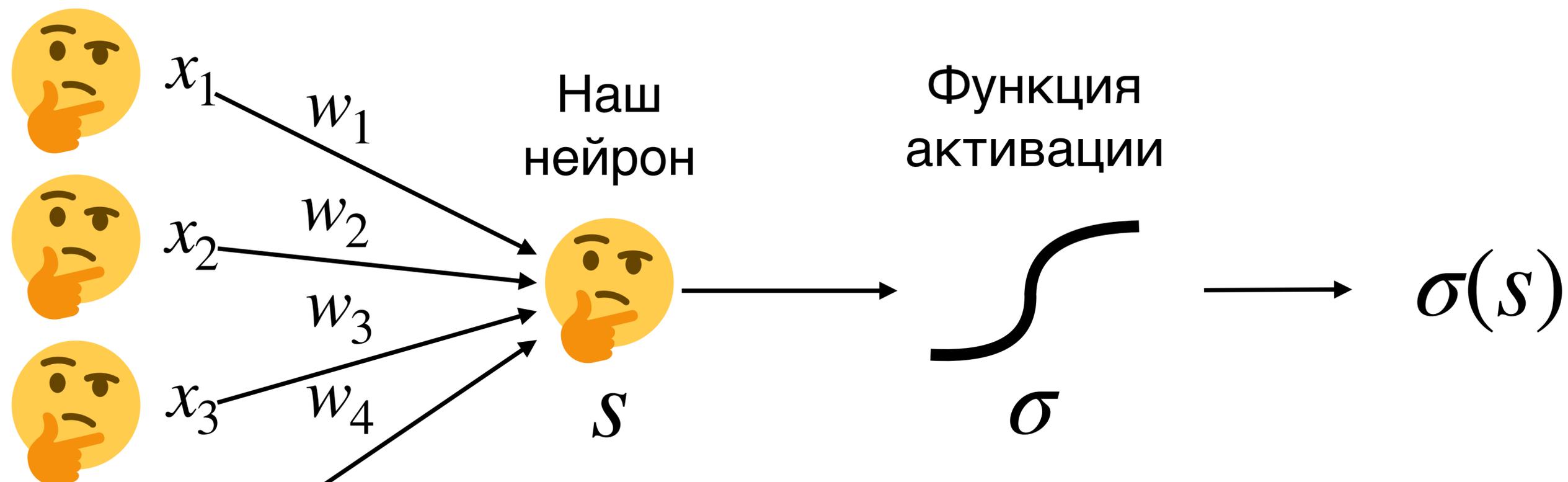
# Как работает 🤖🧠 ?

Другие  
нейроны



# Как работает 🤖🧠 ?

Другие  
нейроны



$$s = w_1x_1 + \dots + w_nx_n$$

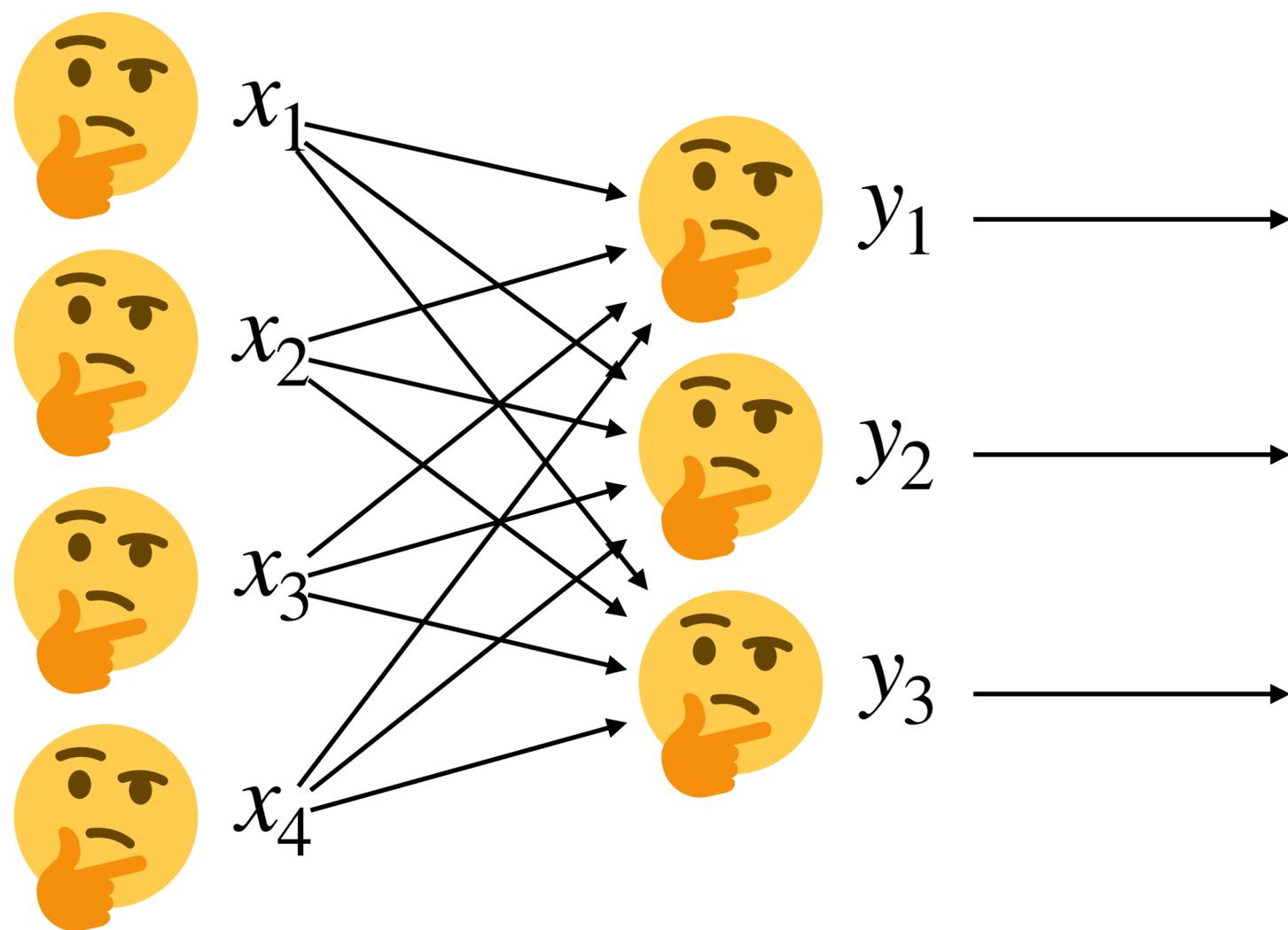
**Вход:**  $x_1$   $x_2$   $x_3$   $x_4$

**Выход:**  $\sigma(w_1x_1 + \dots + w_nx_n)$

**Параметры:**  $w_1$   $w_2$   $w_3$   $w_4$

# Как работает 🤖🧠 ?

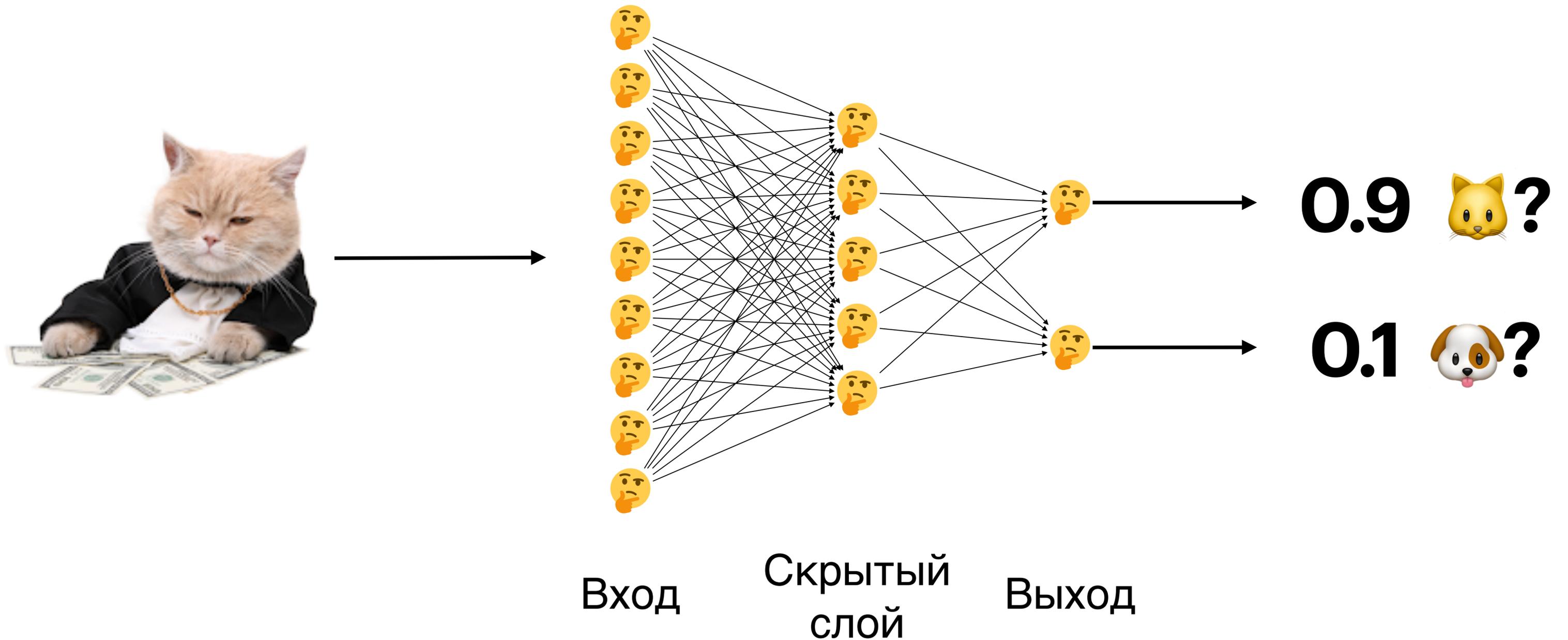
Возьмем несколько нейронов и получим полносвязный слой



$$y = \sigma(Wx)$$

3x1                      3x4 4x1

# Как использовать 🤖🧠 ?



# Обучение нейросети

Объект  
 $x \in \mathbb{R}^d; y \in \mathbb{R}^c$



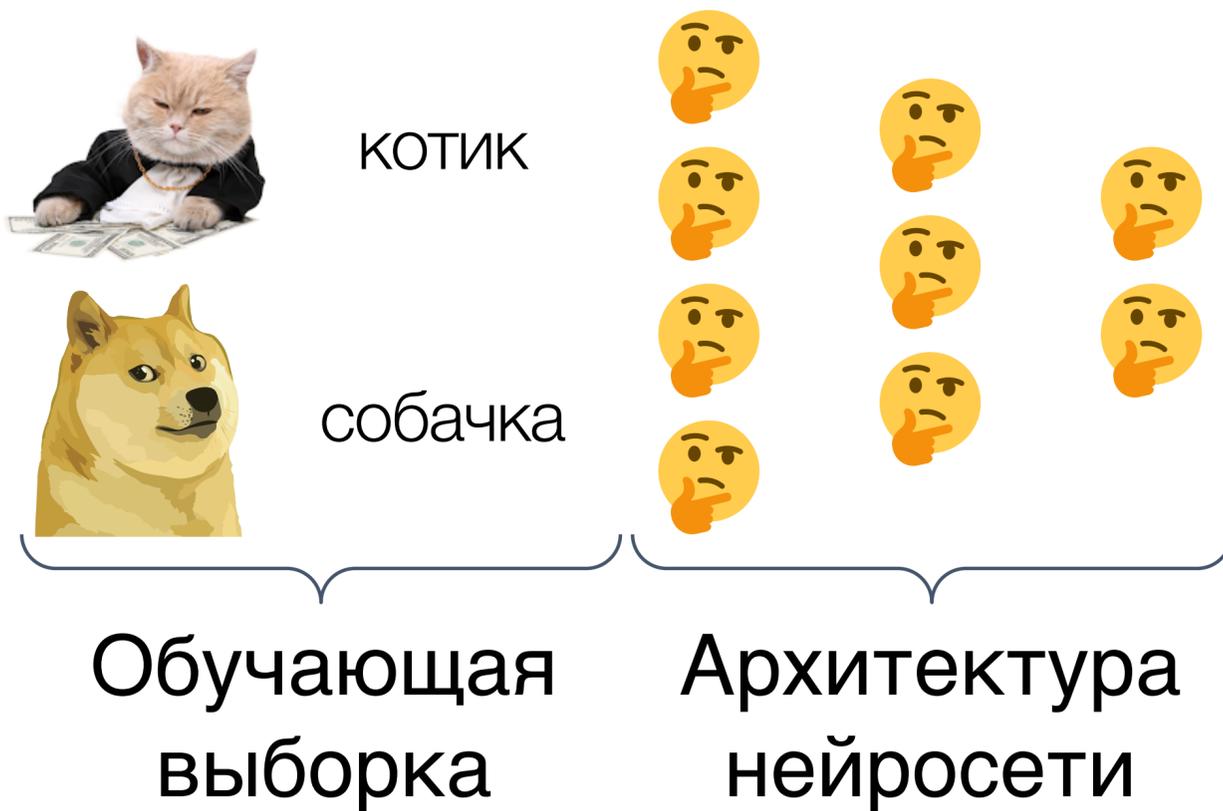
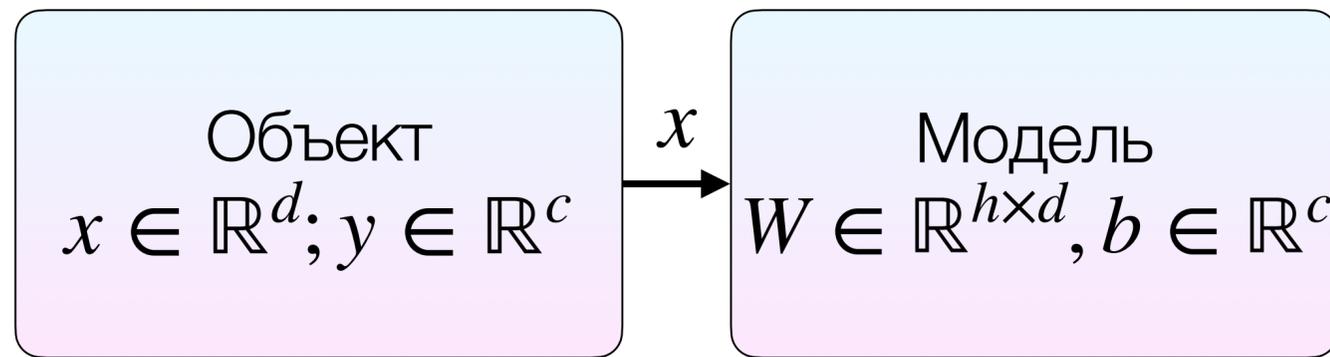
КОТИК



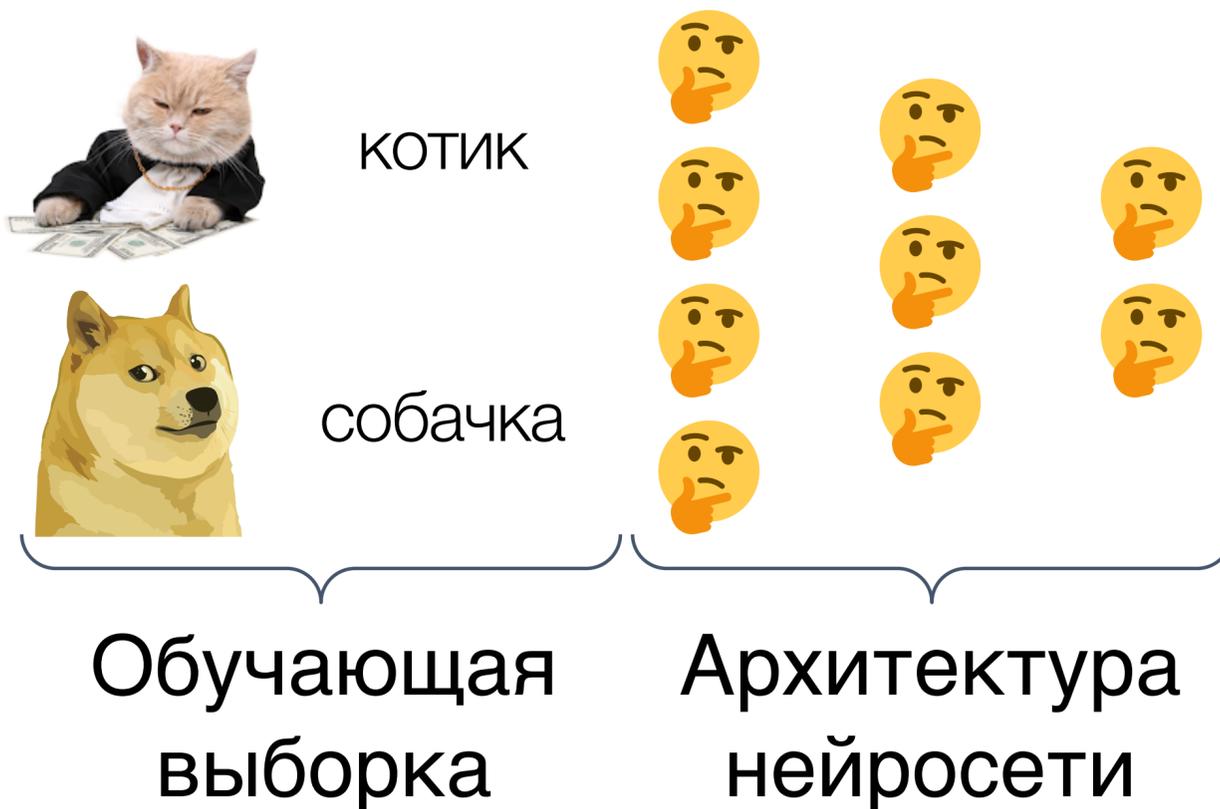
собачка

Обучающая  
выборка

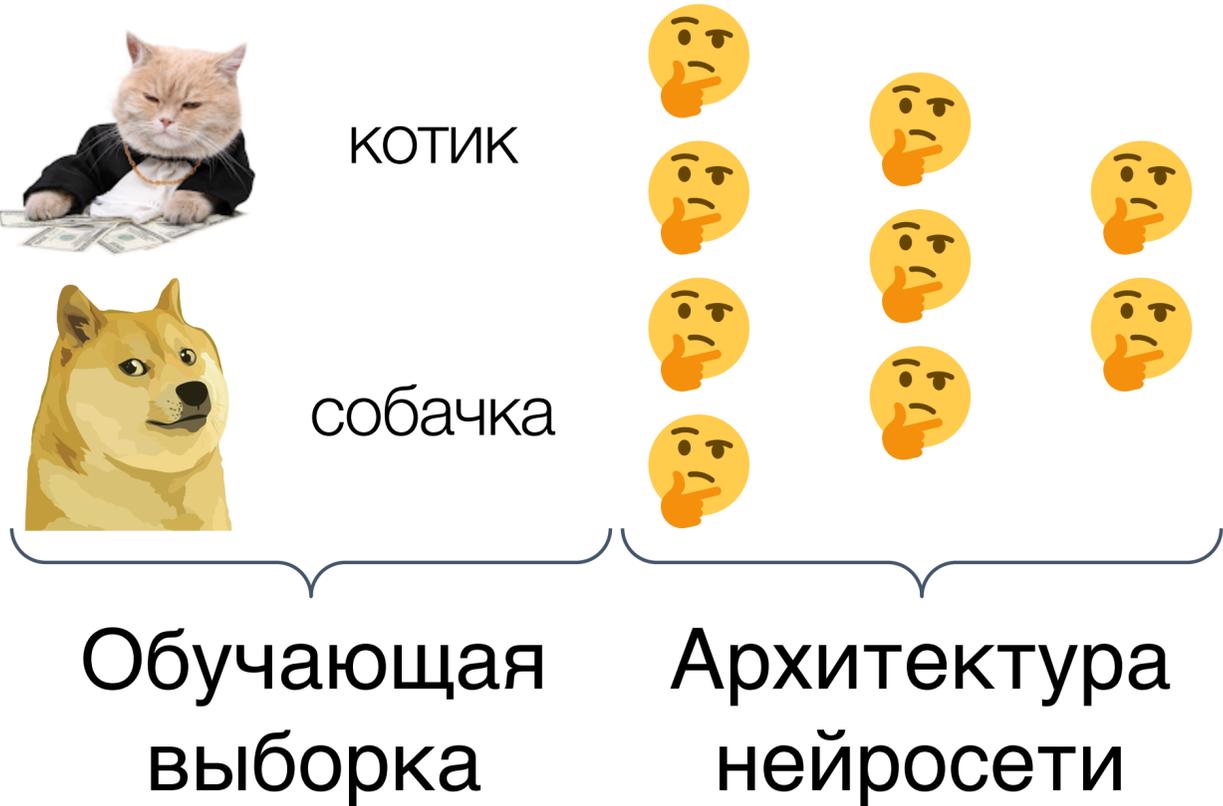
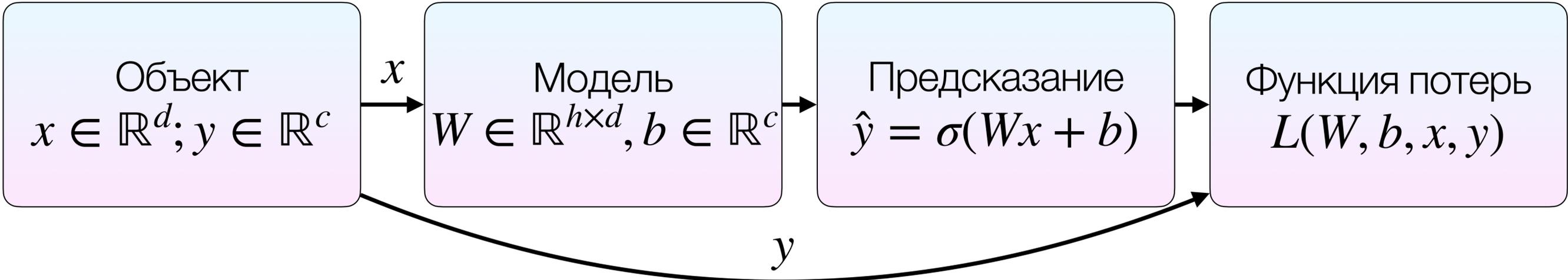
# Обучение нейросети



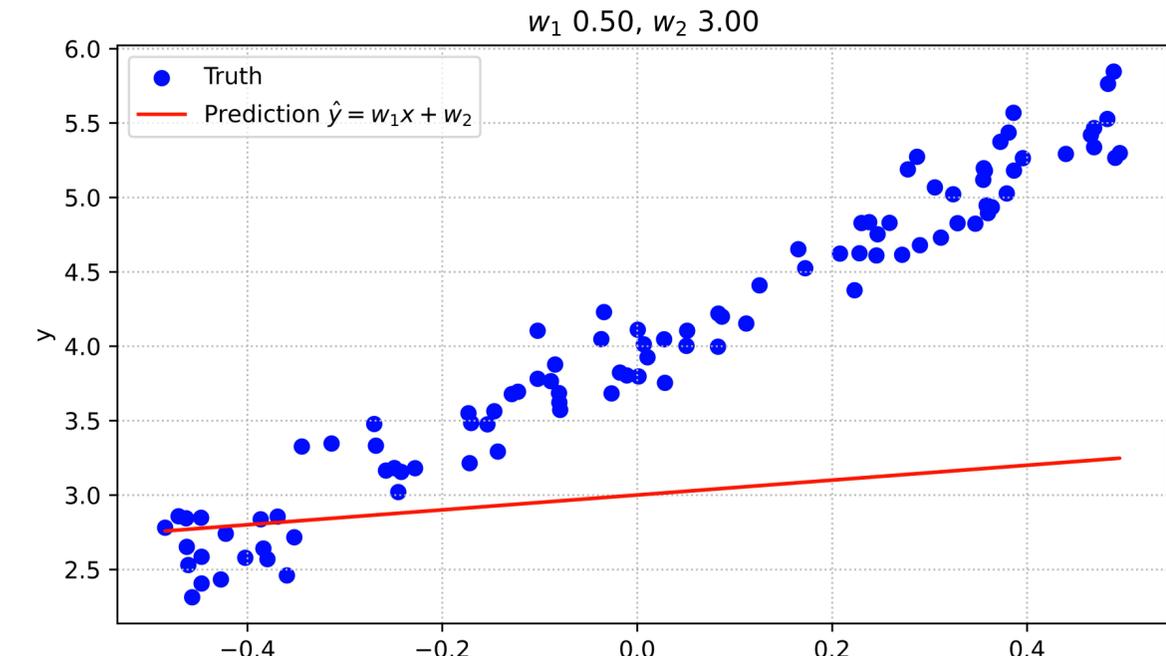
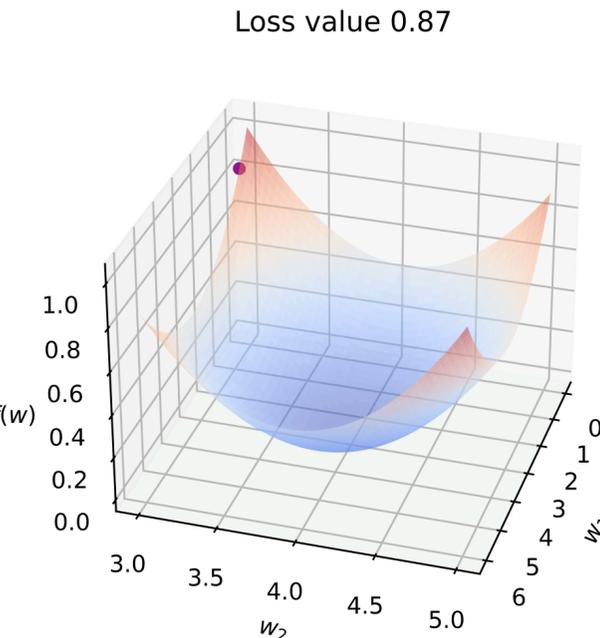
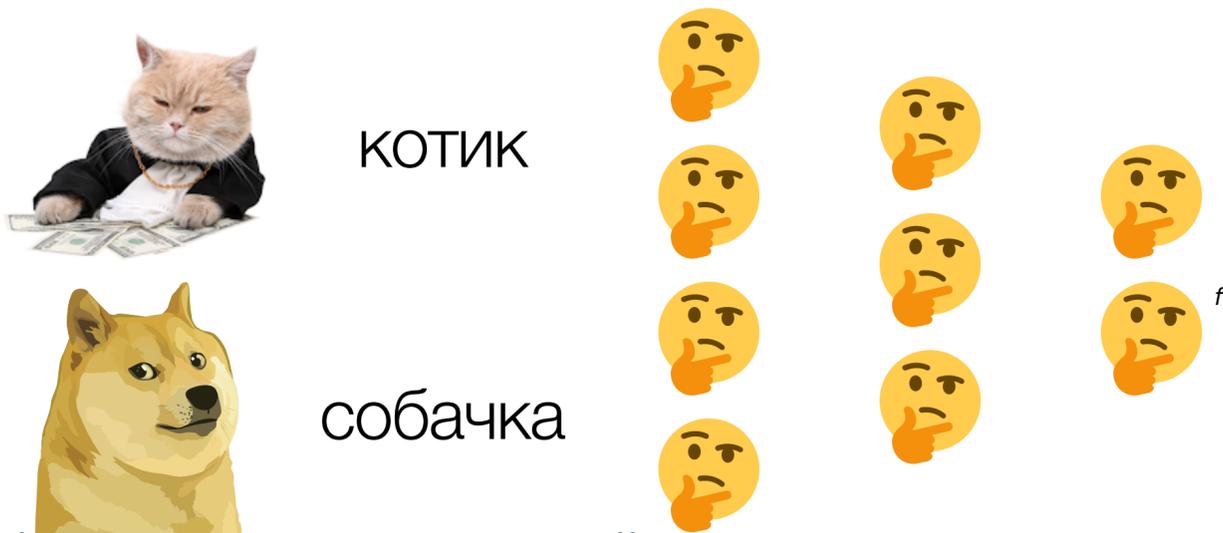
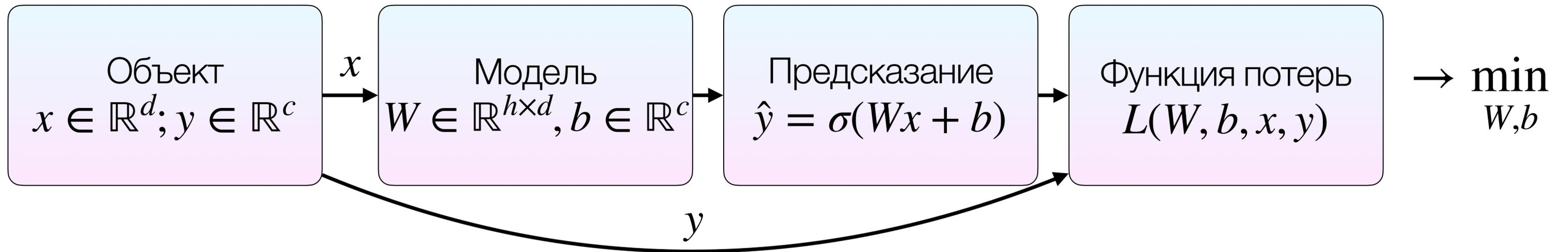
# Обучение нейросети



# Обучение нейросети

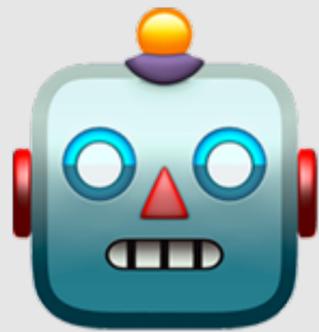


# Обучение нейросети



$$W_{k+1} = W_k - \alpha \nabla_W L(W_k)$$

Изменение параметров таким образом, чтобы уменьшать значение функции потерь на обучающей выборке



**ПРАКТИКА**

# Нехватка памяти для обучения больших моделей

```
RuntimeError: cuda runtime error (2) : out of memory at /data/users/soumith/miniconda2/cond
```

how can i solve this error?



apaszke commented on Mar 8, 2017

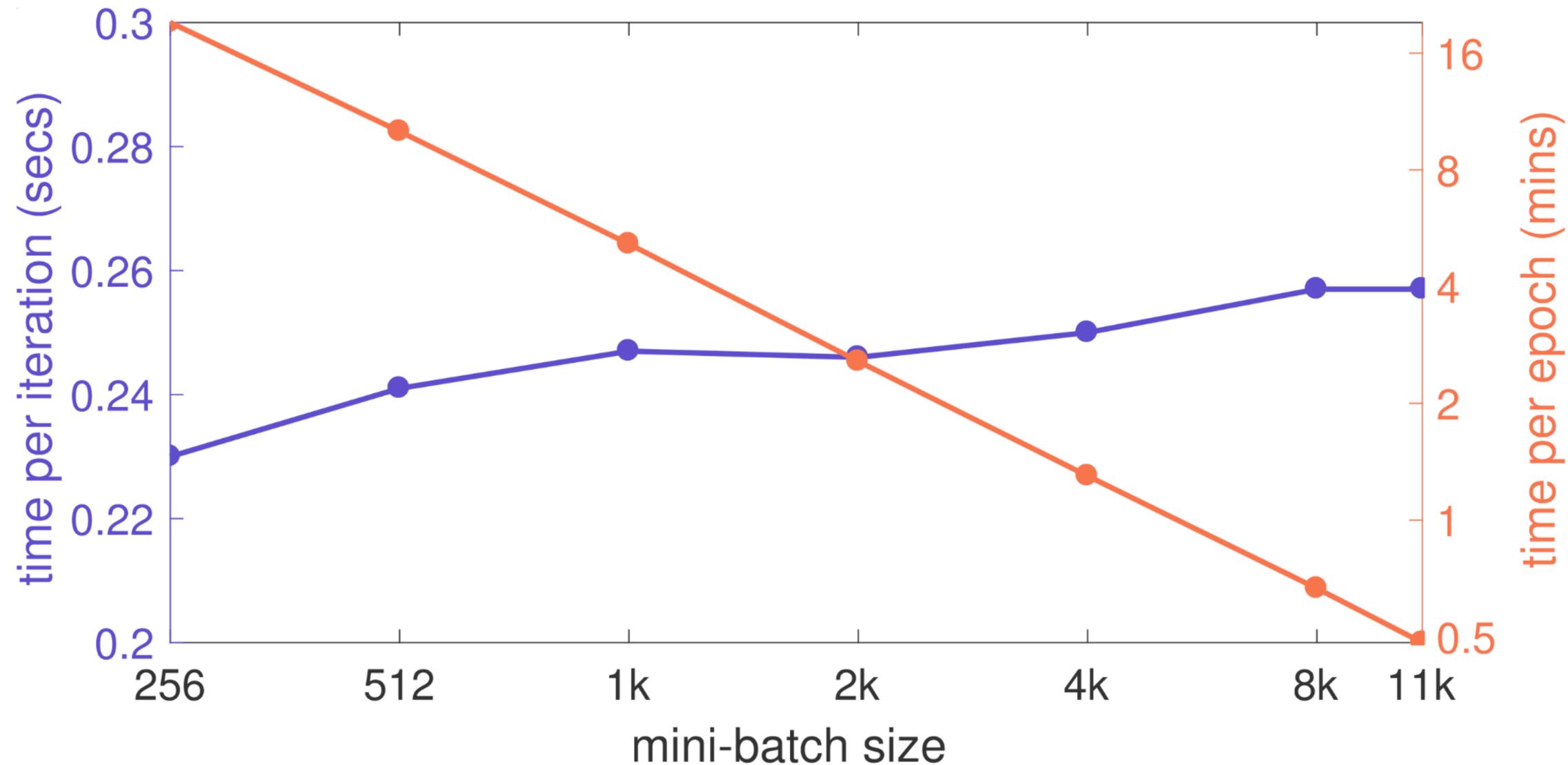
Member



You're running out of memory on the GPU. It's not a bug.



# Размер батча и время, затрачиваемое на одну эпоху



При наличии достаточной памяти GPU, увеличение размера батча позволяет утилизировать ресурсы параллельных вычислений.

 Paper

Accurate, Large Minibatch SGD: Training ImageNet in 1 Hour.

# Просто так увеличить размер батча не получится

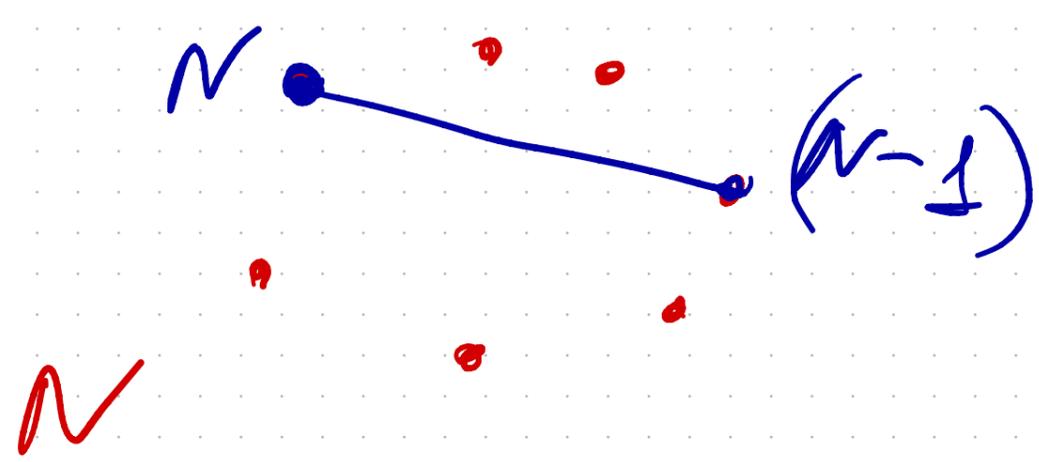
$kn$	$\eta$	top-1 error (%)
256	0.05	23.92 $\pm$ 0.10
256	0.10	23.60 $\pm$ 0.12
256	0.20	23.68 $\pm$ 0.09
8k	0.05 $\cdot$ 32	24.27 $\pm$ 0.08
8k	0.10 $\cdot$ 32	23.74 $\pm$ 0.09
8k	0.20 $\cdot$ 32	24.05 $\pm$ 0.18
8k	0.10	41.67 $\pm$ 0.10
8k	0.10 $\cdot$ $\sqrt{32}$	26.22 $\pm$ 0.03

Обучение ResNet-50 на датасете ImageNet с разными вариантами увеличения размера батча.

(a) **Comparison of learning rate scaling rules.** A reference learning rate of  $\eta = 0.1$  works best for  $kn = 256$  (23.68% error). The linear scaling rule suggests  $\eta = 0.1 \cdot 32$  when  $kn = 8k$ , which again gives best performance (23.74% error). Other ways of scaling  $\eta$  give worse results.

 Paper

Accurate, Large Minibatch SGD: Training ImageNet in 1 Hour.



$$\frac{N \cdot (N-1)}{2}$$

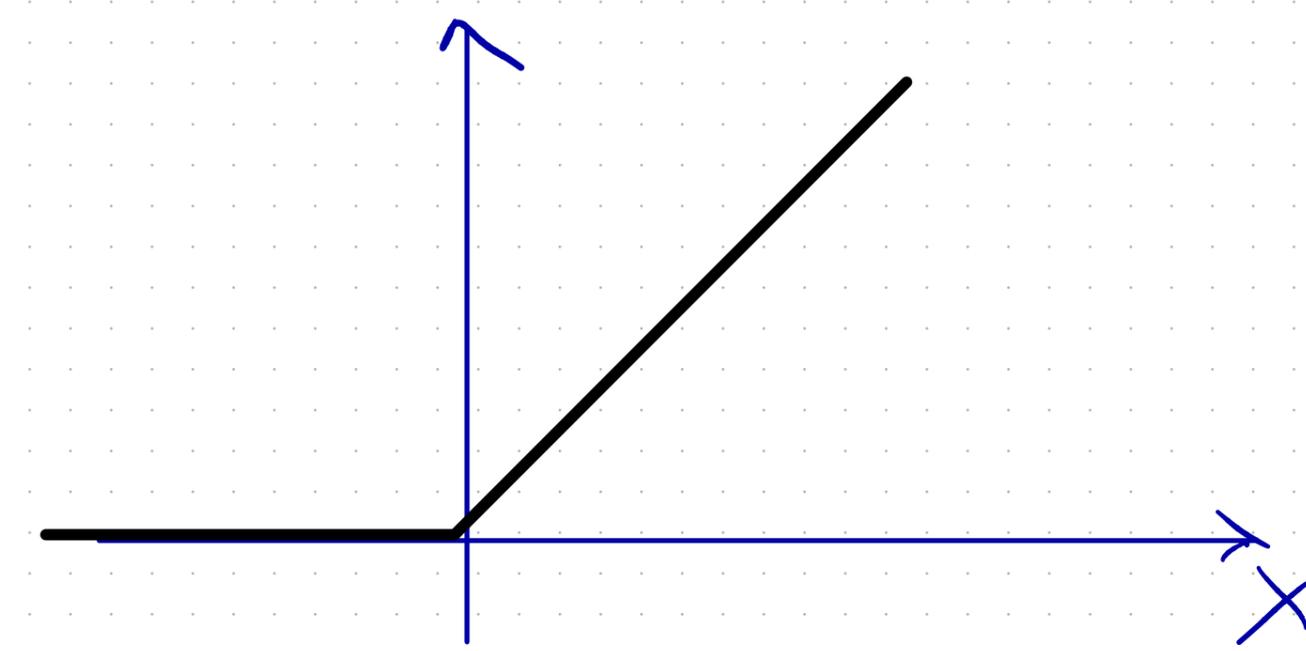
способов

связей

$10^{10}$  нейронов

$$\frac{10^{10} \cdot 10^{10}}{2} \approx 10^{20}$$

ReLU  
Rectified Linear Unit



НЕ ЛИНЕЙНАА

100 слоев

Выход

• взрыв / затухание градиентов

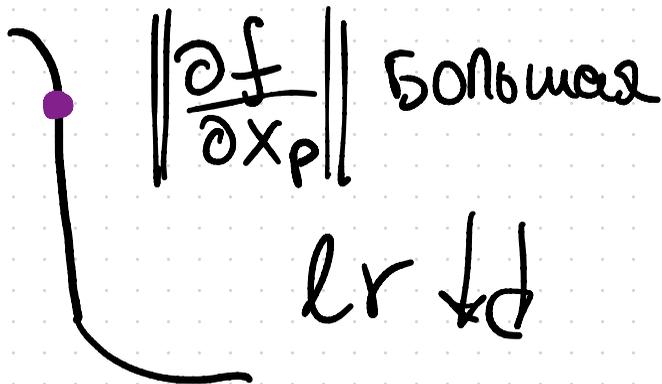
$$\sigma_{100} \left( W_{100} \cdot \sigma_{99} \left( W_{99} \cdot \dots \cdot \sigma_1 (W_1 X) \right) \right)$$

$W_{100} \cdot W_{99} \cdot \dots \cdot W_1 \cdot X$

$$x_{k+1} = x_k - [\nabla^2 f(x_k)]^{-1} \cdot \nabla f(x_k)$$

Пусть  $\nabla^2 f(x_k) = \text{diag}(m_k) = \begin{pmatrix} m_k^1 & & 0 \\ & \dots & \\ 0 & & m_k^n \end{pmatrix}$

$$[\nabla^2 f(x_k)]^{-1} = \text{diag}\left(\frac{1}{m_k}\right)$$



$\left\| \frac{\partial f}{\partial x_p} \right\|$  маленькая



① Для прототипирования имеет смысл выбрать любой адаптивный алгоритм (AdamW, Nadam, Adam, RMSprop...) с дефолтными гиперпараметрами.

② Тем не менее, часто хорошо настроенный SGD + моментум показывает лучшие результаты по сравнению с тюненными Adam / AdamW.

Проекция функции потерь нейронной сети

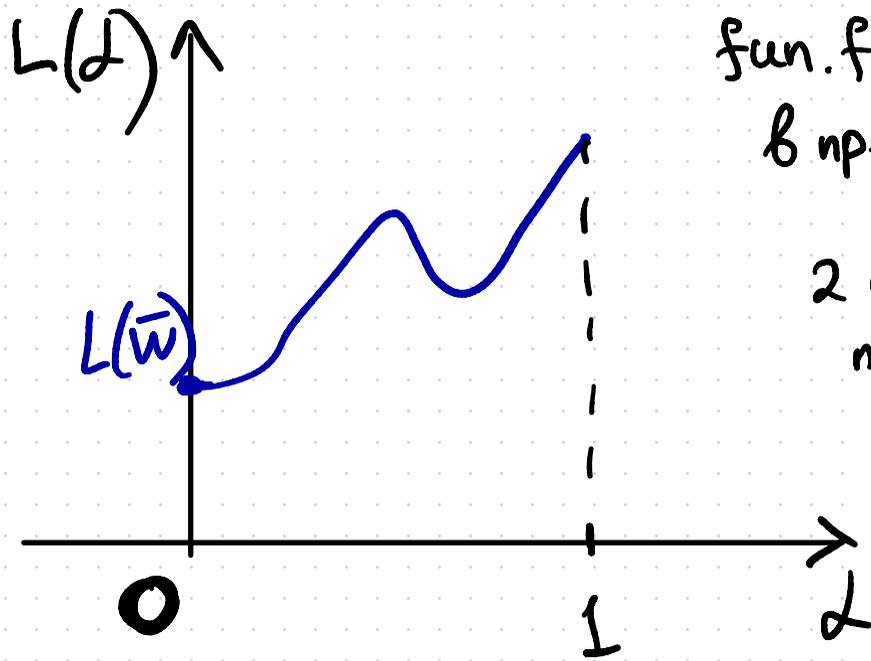
$$L(w) \rightarrow \min_{w \in \mathbb{R}^p}$$

$\bar{w}$  получился после 10 эпох SGD + momentum.

$$w(\alpha) = \bar{w} + \alpha \cdot w_{\perp}$$

можно посчитать  $w(\alpha)$  для  $w_{\perp}$  - случайный вектор  $\in \mathbb{R}^p$   
 $\alpha = \text{np.linspace}(0, 1)$

$$L(w(\alpha))$$



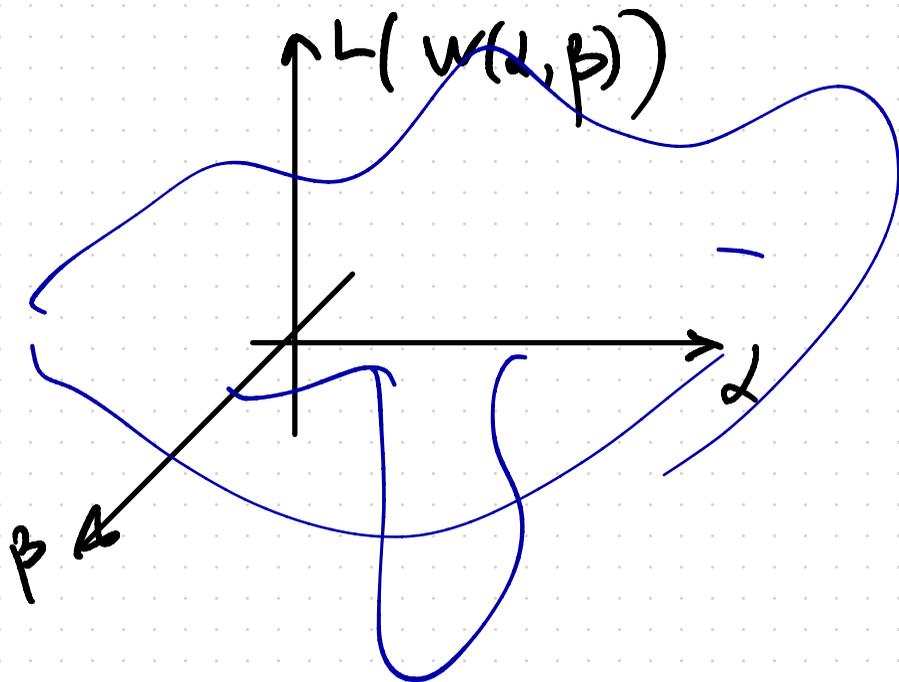
fun. fact

в пр-ве большой  
размерности

2 слуг. вектора  
почти ортогональны

$$w(\alpha, \beta) = \bar{w} + \alpha \cdot w_1 + \beta \cdot w_2$$

$w_1, w_2$  - слуг.  
вектора



Generalization  
error  
SHARPNESS  
AWARE  
MINIMIZATION  
(SAM)

если всё пр-во  $\begin{matrix} 1 \rightarrow \\ 2 \rightarrow \\ 3 \rightarrow \end{matrix}$   $\begin{matrix} 5000 \cdot 6 & z_1 \\ 4000 \cdot 6 & z_2 \\ 2000 \cdot 6 & z_3 \end{matrix}$

Если фабрика за неделю производит:

$$y_1 + y_2 + y_3 \leq \frac{z_1 + z_2 + z_3}{3}$$

фабрика производит в неделю:

capacity

$$\frac{y_1}{z_1} + \frac{y_2}{z_2} + \frac{y_3}{z_3} \leq 1$$

