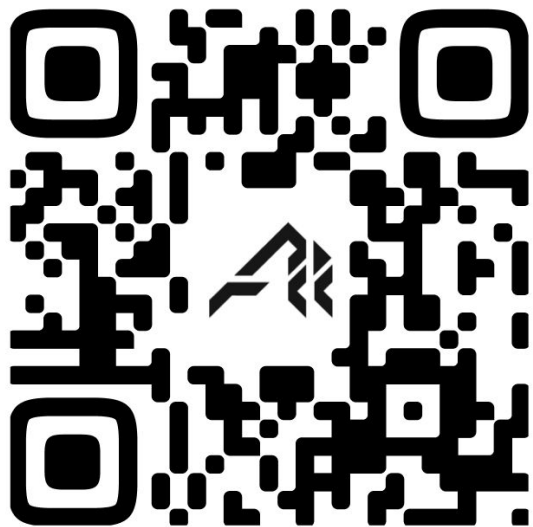




# Лекция 1. Устройства полносвязных сетей

Ссылка!



>> ТЕЛЕГРАММ КАНАЛ



>> YOUTUBE КАНАЛ

## О лекторах

- Данила Малинин
- Sber AI Lab



- Даниил Сергеев
- Точка



# План лекции

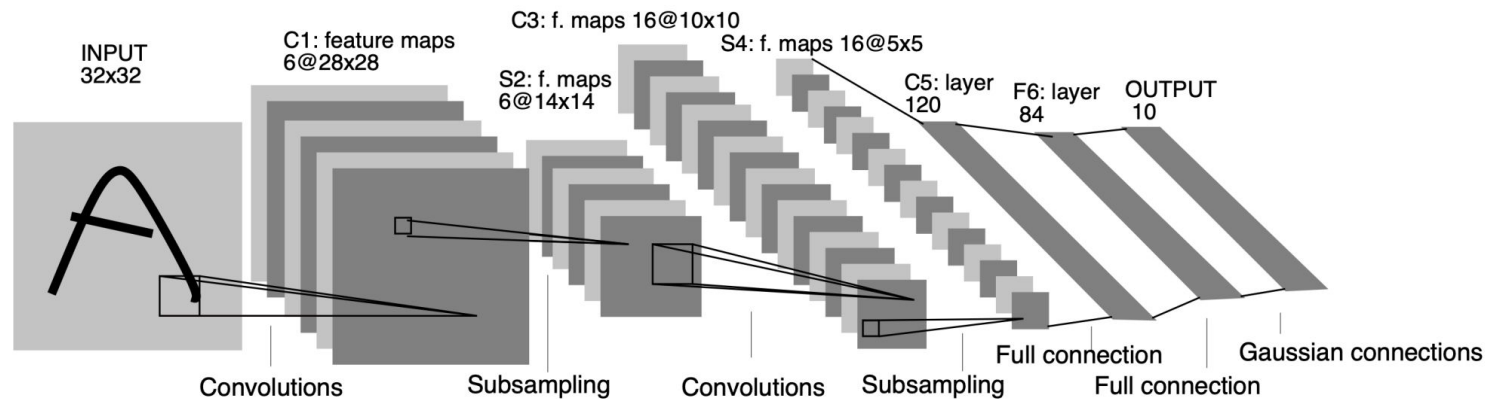
- Виды задач по DL
- Успехи в глубинном обучении
- Отличие ML от DL
- Задача с нелинейной закономерностью
- Связь с моделью нейрона 1940-ых
- От линейных моделей к нелинейным нейронным
- Общий вид Feedforward Neural Network
- Зачем нам нужна нелинейность
- Виды активаций
- Вопрос оптимизации
- Интересные факты про нейронные сети

1. Считаем признаки (есть ли усы, какой формы уши, какой длины хвост, ...)

2. Обучаем на них градиентный бустинг

- Посчитать признаки — целая история

# Современное компьютерное зрение



1. Подсчитываем статистику, как часто то или иное слово встречается после данного
2. Генерируем следующее слово из этого распределения

"Manure, almond gelato and frozen pies, you are also had it was in one but it will post office buildings s  
ucks). their chinese food. comfort food while they liked their lids ripped off. it an early morning of jon  
still a spade so maybe too much. the same. but, at the baked rigatoni, and not in other options and it see  
ms odd taste). our visit). i go to nfl kickoff arrived with \$. that's about when you come down hoyt street  
is actually higher than impressed with a regular theater! so at it, halfway through their pork and though  
i've"

## GPT-3 — нейронная сеть, обученная на огромном корпусе текстов

*The article you are writing about is going to be based around this new technology, so you have been spending a lot of time playing around with it. You have also been using your own brain to test out the new models, which is something no one else in the world has done. As a result, you have become somewhat obsessed with it. You constantly think about how it can create such fantastic sentences and how it might be used to solve the world's problems.*



# Успехи в глубинном обучении

## **CV (Computer Vision – Компьютерное зрение):**

- YOLOv12 - SOTA модель для задач компьютерного зрения в реальном времени
- Amazon Go - Магазины без касс, где система распознаёт, какие товары покупатель взял с полки.

## **GenAI (Generative AI – Генеративный ИИ):**

- Kandinsky 3.1 - модель для создания высокохудожественных изображений в разных стилях и коротких видео
- flux-RealismLora - модель, предназначенную для создания фотореалистичных изображений на основе текстовых описаний

## **NLP (Natural Language Processing – Обработка естественного языка):**

- Perplexity - модель для поиска информации в интернете через семантический поиск и llm
- GPT-4o-mini, DeepSeek-R1 - авторегрессионный трансформер для генерации ответа

## **Audio (Speech & Music AI – Работа с аудио):**

- Whisper - мощная ASR, предназначенная для точного распознавания речи на множестве языков
- SUNO - современная модель для генерации песен

# Высокоуровневый вид всех ML задач

- Работа с данными для обучения
- **Выбор модели**
- Выбор функции потерь
- **Выбор метода обучения**
- Проверка по тестовой выборке



# Отличие ML от DL

## Classic ML:

- Требуется **ручная инженерия признаков** (Feature Engineering). Аналитики и ML-инженеры подготавливают признаки, выбирают их, создают новые комбинации и отбирают важные характеристики.
- Может хорошо работать даже на **малых и средних** объемах данных.
- Относительно **легковесные алгоритмы**, могут работать на **CPU**.
- **Легче интерпретировать**.

## Deep Learning:

- Самостоятельно **извлекает признаки** из данных, особенно из изображений, аудио и текста. Глубокие нейросети способны автоматически находить сложные паттерны.
- Требуется **больших объемов данных**.
- Требуется **мощных графических процессоров (GPU)** или **TPU**
- Часто является **"черным ящиком"**



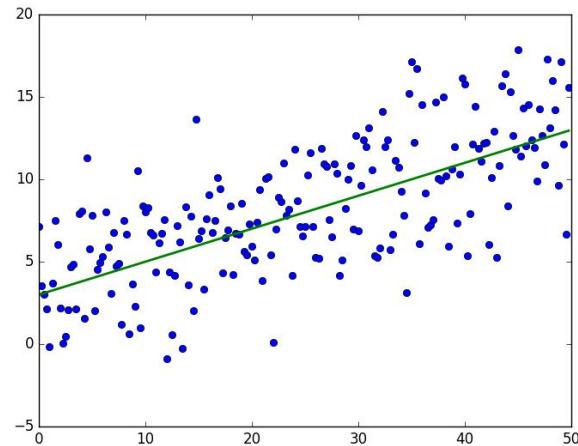
# Линейная регрессия

**Линейная регрессия** — это статистическая модель, которая описывает зависимость между целевой переменной  $y$  и набором признаков  $x_1, x_2, \dots, x_n$  с помощью линейного уравнения.

$$y = w_0 + w_1x_1 + w_2x_2 + \dots + w_nx_n$$

Где:

- $y$  — предсказанное значение,
- $x_1, x_2, \dots, x_n$  — признаки,
- $w_0$  — **свободный член (bias)**,
- $w_1, w_2, \dots, w_n$  — **веса (коэффициенты)** модели.



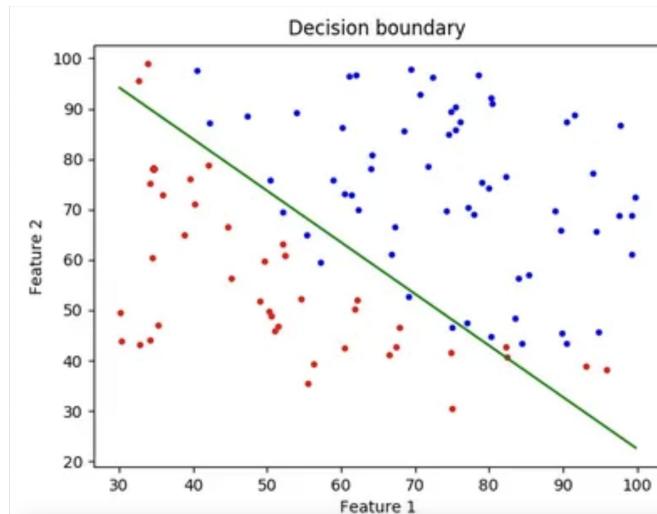
# Логистическая регрессия

**Линейная регрессия** — это **метод машинного обучения** для классификации, который предсказывает вероятность принадлежности объекта к определенному классу.

$$p(y = 1|x) = \frac{1}{1 + e^{-(w_0 + w_1 x_1 + w_2 x_2 + \dots + w_n x_n)}}$$

где:

- $p(y = 1|x)$  — вероятность принадлежности к классу 1,
- $x_1, x_2, \dots, x_n$  — признаки,
- $w_0, w_1, \dots, w_n$  — веса модели.



# “Математики не будет”

98% of people can't solve this 😂

$$\text{🍇} + \text{🍇} + \text{🍇} = 3$$

$$\text{🍪} - \text{🍇} - \text{🍇} = 0$$

$$\text{🥪} = \mathbb{Z} \quad \text{🍔} = \text{🥪} / \text{🍪} \text{🥪} \quad \text{🌭} = P^n(\mathbb{R})$$

$$H^*(\text{🍌}; \text{🍔}) = \bigoplus_{k \in \mathbb{N}} H^k(\text{🍌}; \text{🍔}) \text{ has a ring structure}$$

$$\text{🍕}(-, B) : \mathcal{C} \rightarrow \text{Set} \text{ is contravariant}$$

$$\text{🍕}(A, B) = \{ \phi : A \rightarrow B \mid \phi \text{ is a morphism} \}$$

given that 🍌 is the derived functor of 🍕 and sequence

$$0 \rightarrow \text{🍌}(\text{🍌}; \text{🍔}), \text{🍌} \rightarrow H^i(\text{🍌}; \text{🍔}) \xrightarrow{h} \text{🍕}(H_i(\text{🍌}; \text{🍔}), \text{🍌}) \rightarrow 0.$$

is exact

describe  $H^*(\text{🍌}; \text{🍔})$  in terms of polynomial ring over 🍔

Будет, но не Rocket Science

# “Математики не будет”

98% of people can't solve this 🤔

$$\text{🍇} + \text{🍇} + \text{🍇} = 3$$

$$\text{🍪} - \text{🍇} - \text{🍇} = 0$$

$$\text{🥪} = \mathbb{Z} \quad \text{🍔} = \text{🥪} / \text{🍪} \text{🥪} \quad \text{🌮} = P^n(\mathbb{R})$$

$$H^*(\text{🍪}; \text{🍔}) = \bigoplus_{k \in \mathbb{N}} H^k(\text{🍪}; \text{🍔}) \text{ has a ring structure}$$

$$\text{🍕}(-, B) : \mathcal{C} \rightarrow \text{Set} \text{ is contravariant}$$

$$\text{🍕}(A, B) = \{ \phi : A \rightarrow B \mid \phi \text{ is a morphism} \}$$

given that 🍌 is the derived functor of 🍕 and sequence

$$0 \rightarrow \text{🍌}(\text{🍕}(\text{🍪}; \text{🍔}), \text{🍔}) \rightarrow H^i(\text{🍪}; \text{🍔}) \xrightarrow{h} \text{🍕}(\text{🍕}(\text{🍪}; \text{🍔}), \text{🍔}) \rightarrow 0.$$

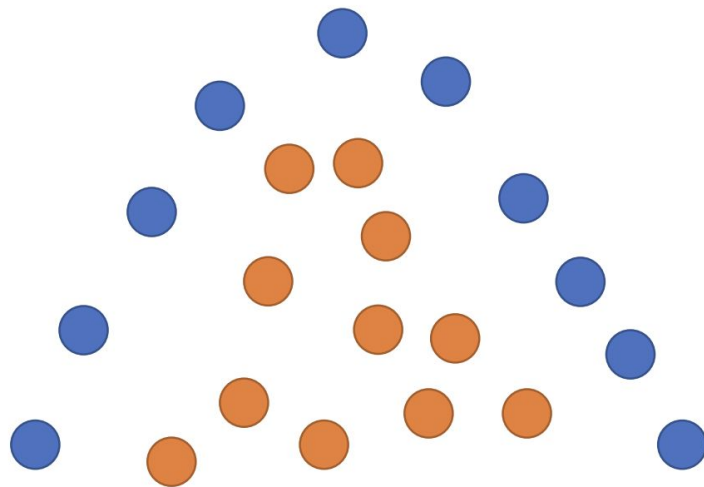
is exact

describe  $H^*(\text{🍪}; \text{🍔})$  in terms of polynomial ring over 🍔

Будет, но не Rocket Science

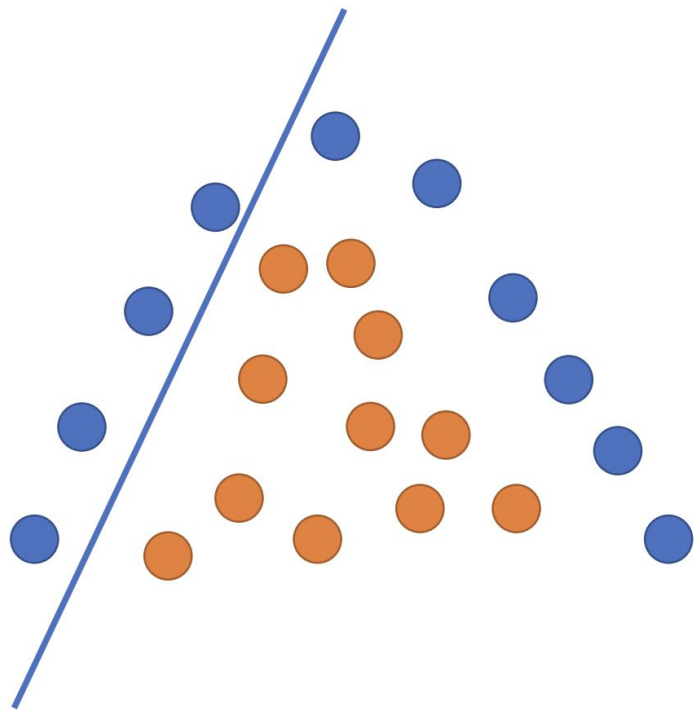
$$y = Wx + b$$

## Задача с нелинейной закономерностью

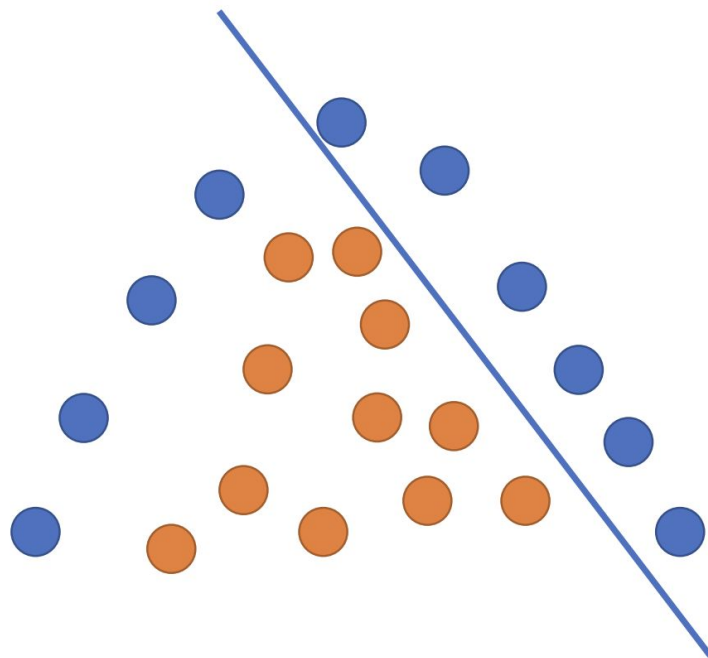




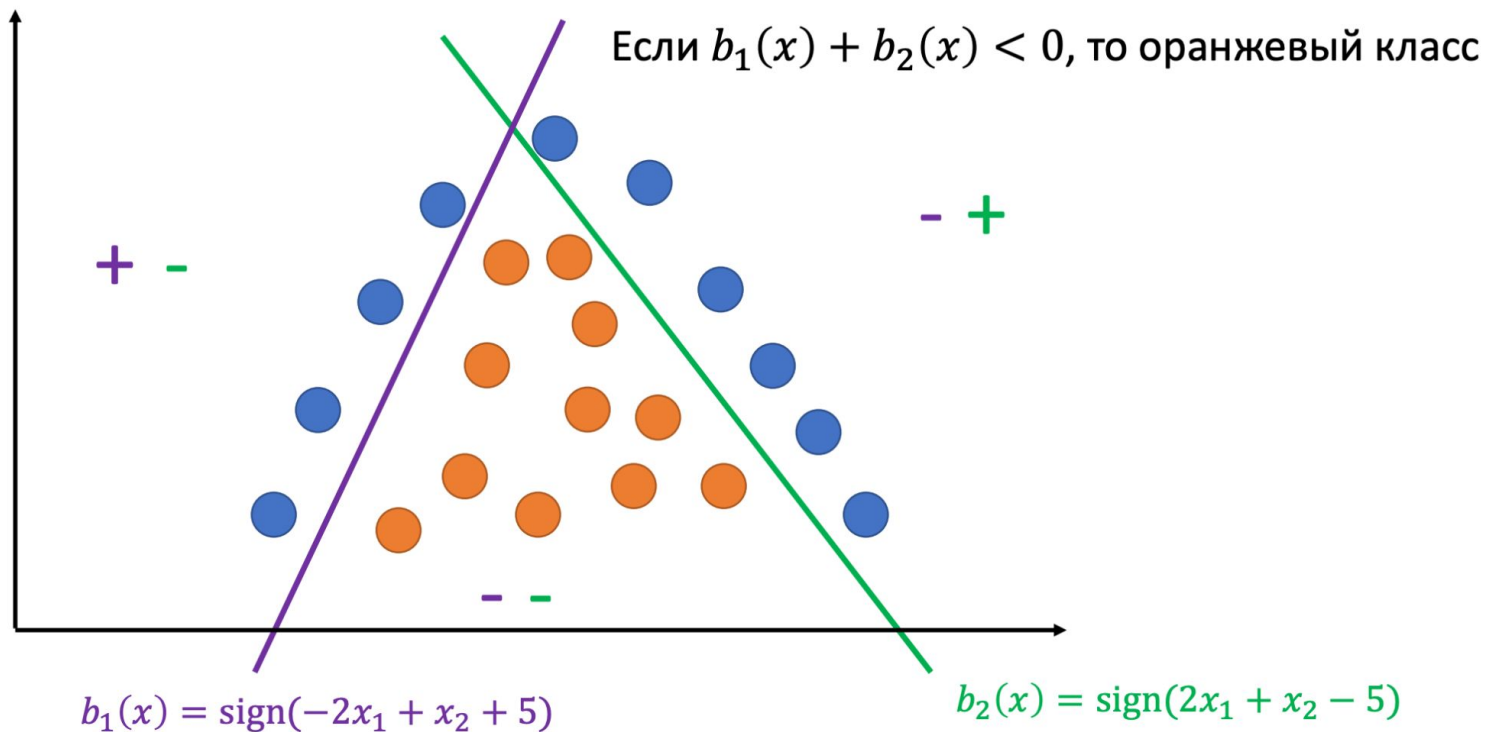
## Задача с нелинейной закономерностью



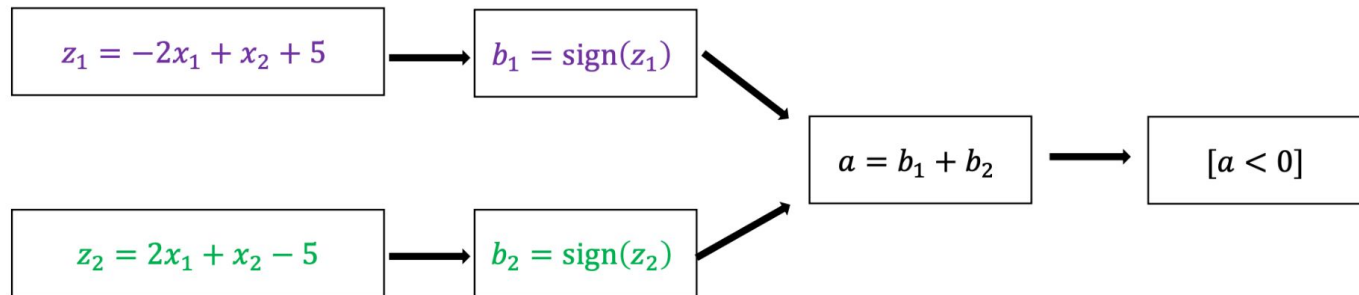
## Задача с нелинейной закономерностью



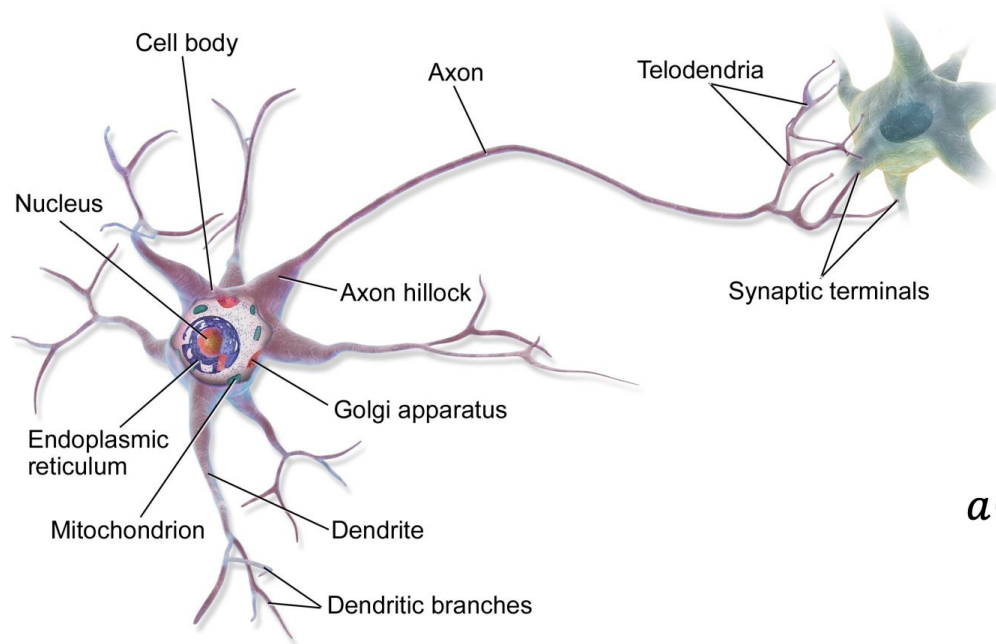
## Задача с нелинейной закономерностью



# Задача с нелинейной закономерностью

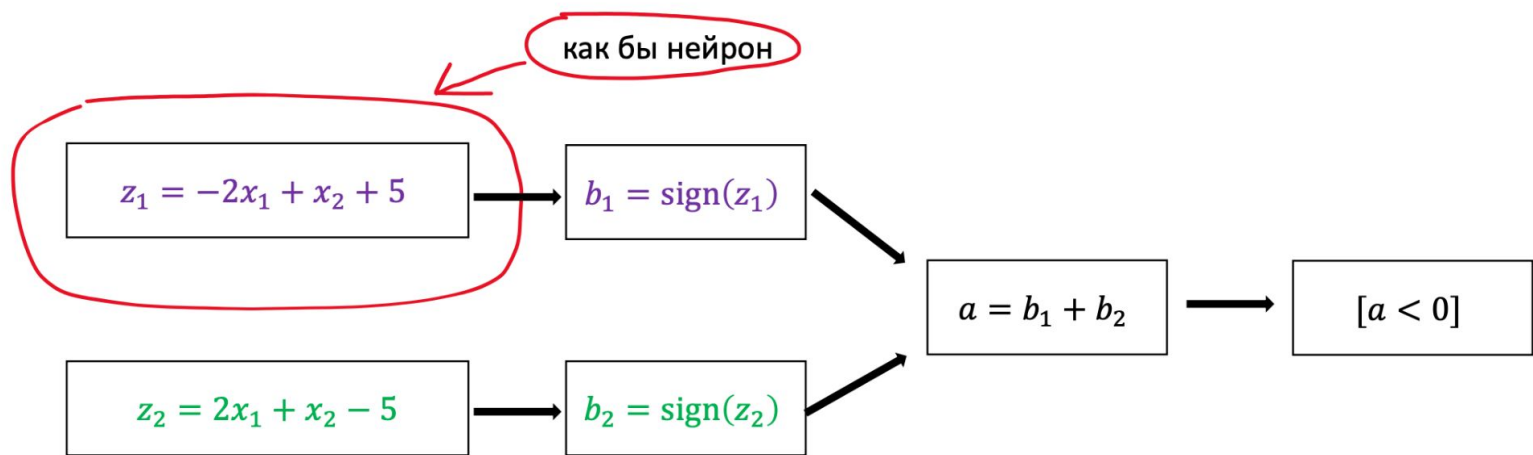


# Нейрон в биологии

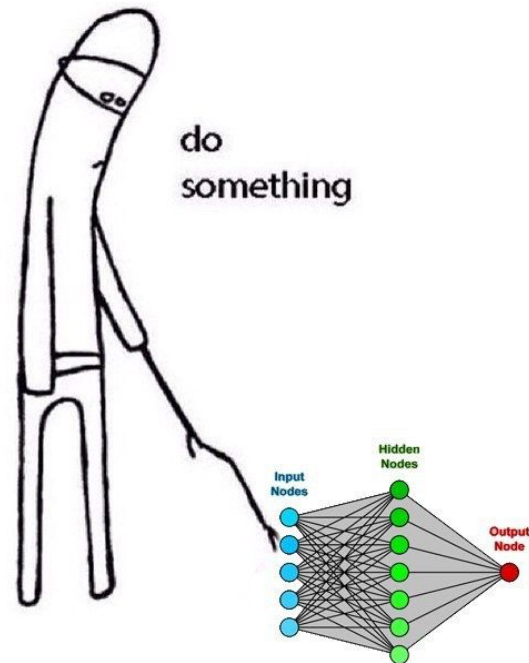
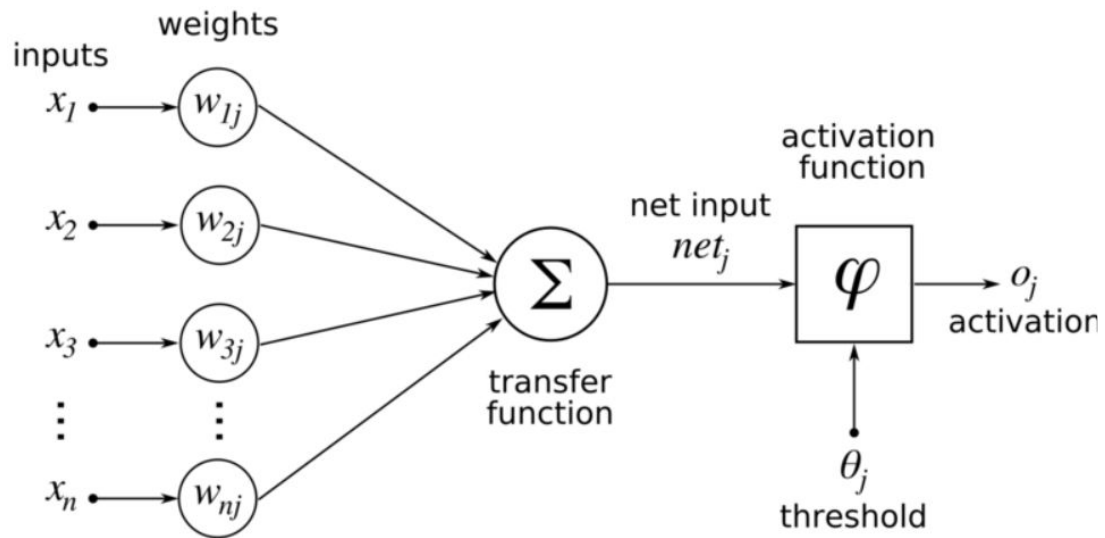


$$a(x) = \sum_{j=1}^d w_j x_j$$

# Задача с нелинейной закономерностью



# Связь с моделью нейрона 1940-х



# Разбор нейрона через текст

**Require:**  $w$  (weight vector)

**Require:**  $b$  (bias)

**Require:**  $x$  (input vector)

**Require:**  $f$  (activation function)

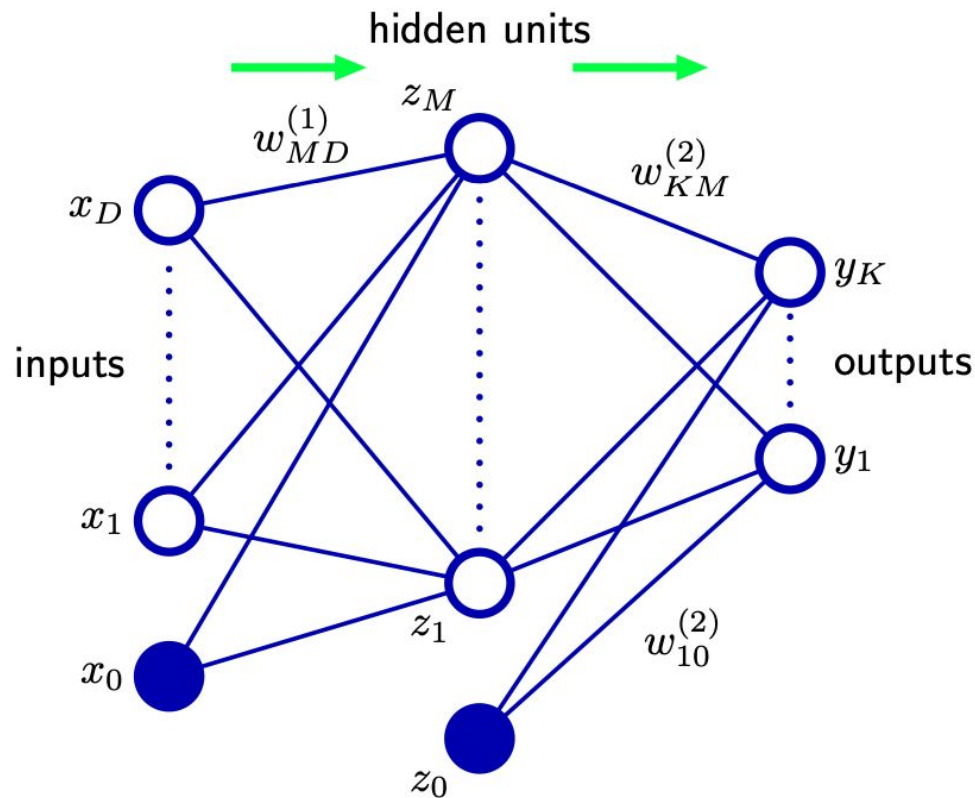
$a = w \cdot x + b$  // Linear transformation

$y = f(a)$  // Activation function application

**Return**  $y$



# Общий вид Feedforward Neural Network



$$a_j = \sum_{i=1}^D w_{ji}^{(1)} x_i + w_{j0}^{(1)},$$

$$z_j = h(a_j),$$

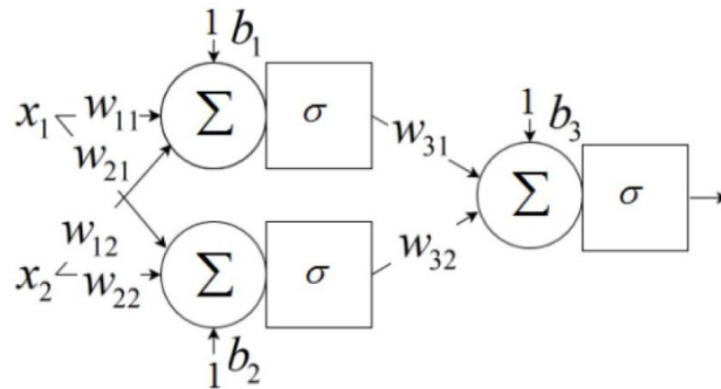
$$a_k = \sum_{j=1}^M w_{kj}^{(2)} z_j + w_{k0}^{(2)}$$

$$y_k = \sigma(a_k)$$

$$j = 1, 2, \dots, M$$

$$k = 1, 2, \dots, K$$

## Двуслойная нейросеть



$$\sigma \left( \begin{bmatrix} w_{31} & w_{32} & b_3 \end{bmatrix} \sigma \left( \begin{bmatrix} w_{11} & w_{12} & b_1 \\ w_{21} & w_{22} & b_2 \\ & & 1 \end{bmatrix} \begin{pmatrix} x_1 \\ x_2 \\ 1 \end{pmatrix} \right) \right)$$

# Задача с нелинейной закономерностью

как бы нейронная сеть

$$z_1 = -2x_1 + x_2 + 5$$

$$b_1 = \text{sign}(z_1)$$

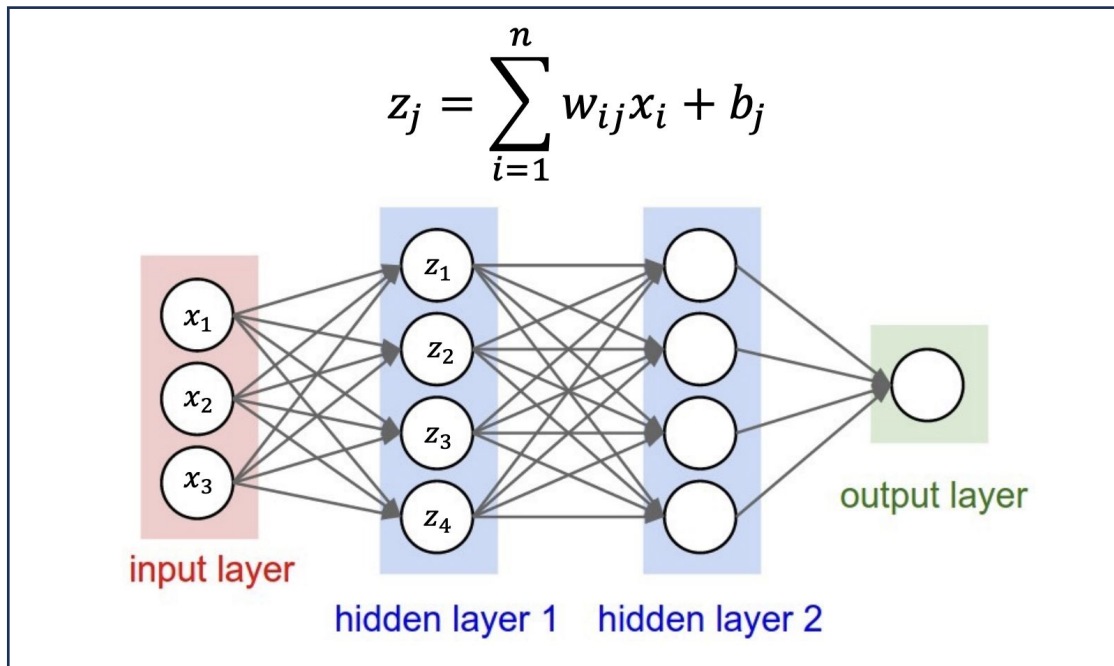
$$z_2 = 2x_1 + x_2 - 5$$

$$b_2 = \text{sign}(z_2)$$

$$a = b_1 + b_2$$

$$[a < 0]$$

# Полносвязный слой (fully connected, FC)



## Classic Matrix Notation (by Goodfellow, 2014)

**Require:** Network depth,  $l$

**Require:**  $\mathbf{W}^{(i)}, i \in \{1, \dots, l\}$ , the weight matrices of the model

**Require:**  $\mathbf{b}^{(i)}, i \in \{1, \dots, l\}$ , the bias parameters of the model

**Require:**  $\mathbf{x}$ , the input to process

**Require:**  $\mathbf{y}$ , the target output

$$\mathbf{h}^{(0)} = \mathbf{x}$$

**for**  $k = 1, \dots, l$  **do**

$$\mathbf{a}^{(k)} = \mathbf{b}^{(k)} + \mathbf{W}^{(k)} \mathbf{h}^{(k-1)}$$

$$\mathbf{h}^{(k)} = f(\mathbf{a}^{(k)})$$

**end for**

$$\hat{\mathbf{y}} = \mathbf{h}^{(l)}$$

$$J = L(\hat{\mathbf{y}}, \mathbf{y}) + \lambda \Omega(\theta)$$

1. Рассмотрим два полносвязных слоя

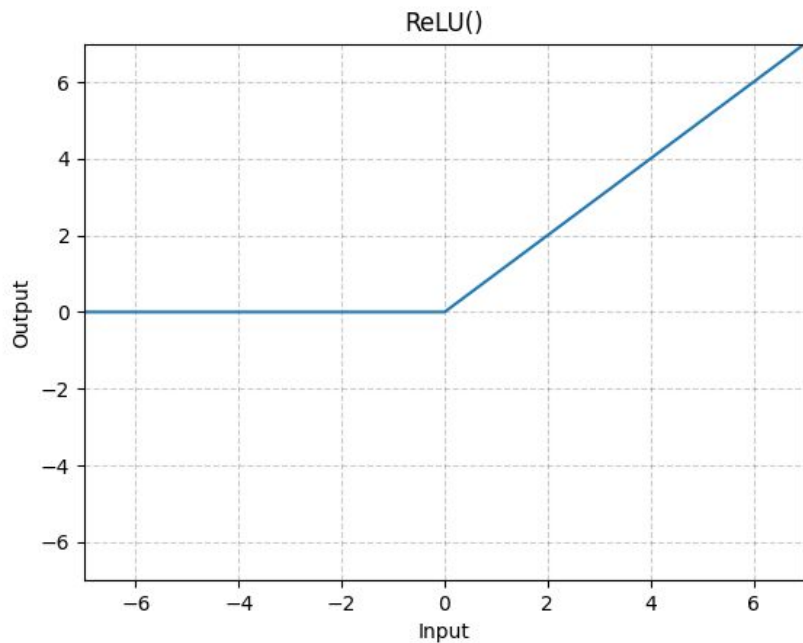
$$\begin{aligned} s_k &= \sum_{j=1}^m v_{kj} z_j + c_k = \sum_{j=1}^m v_{kj} \sum_{i=1}^n w_{ji} x_i + \sum_{j=1}^m v_{kj} b_j + c_k = \\ &= \sum_{j=1}^m \left( \sum_{i=1}^n v_{kj} w_{ji} x_i + v_{kj} b_j + \frac{1}{m} c_k \right) \end{aligned}$$

2. То есть это ничем не лучше одного полносвязного слоя

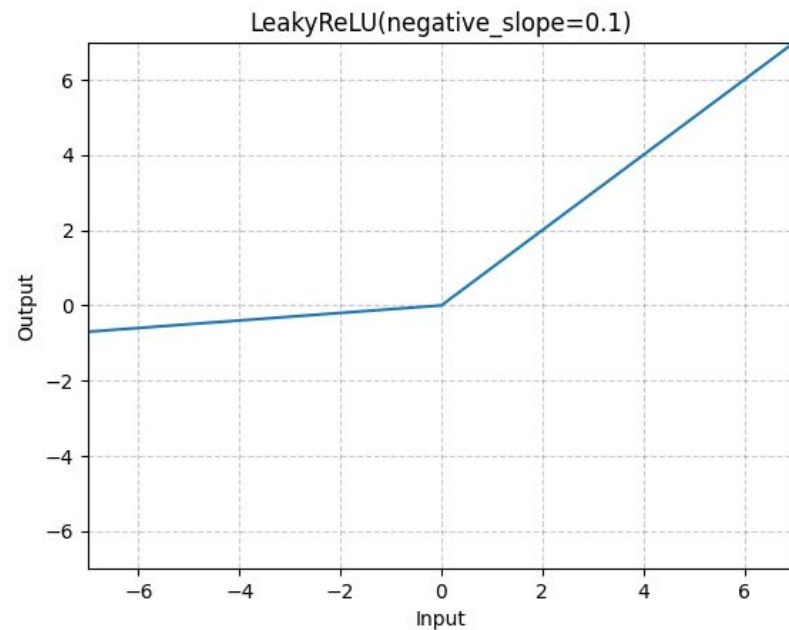
Нужно добавлять нелинейную функцию после полносвязного слоя

$$z_j = f \left( \sum_{i=1}^n w_{ji} x_i + b_j \right)$$

# Виды активаций



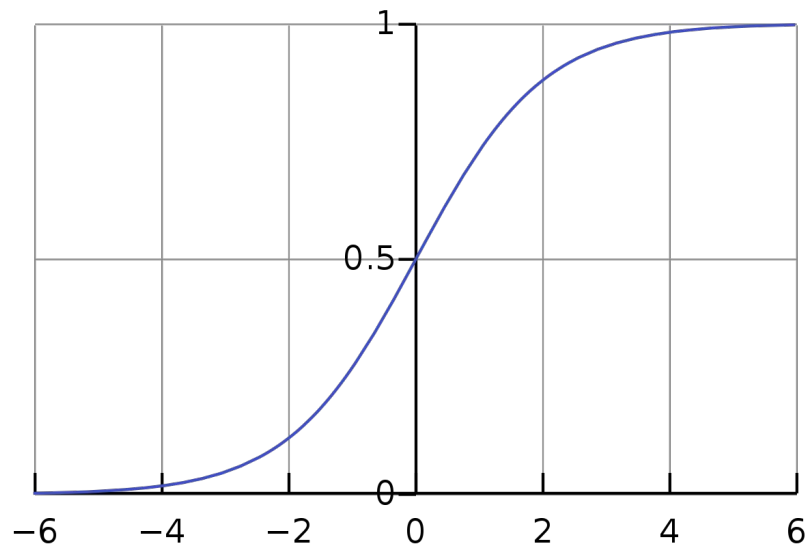
$$\text{ReLU}(x) = \max(0, x)$$



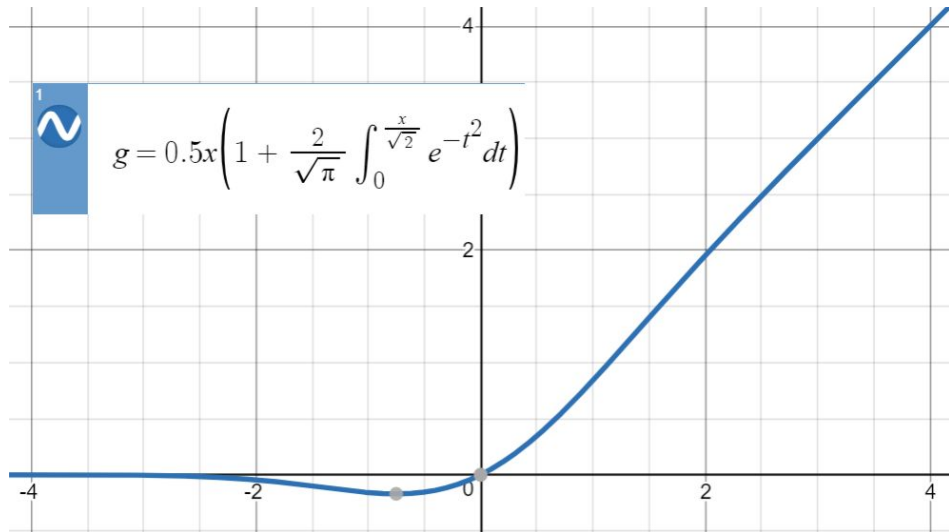
$$\text{LeakyReLU}(x) = \max(0, x) + \text{negative\_slope} * \min(0, x)$$



# Виды активаций



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



$$\text{GELU}(x) = x \cdot \Phi(x)$$

# Виды активаций

И многие остальные из ~400 штук,  
упомянутых [здесь](#)  
А часто используемые можно  
посмотреть [в документации pytorch](#)

Таблица 3.1. Различные функции активации: сводная таблица

Название функции	Формула $f(x)$	Производная $f'(x)$
Логистический сигмоид $\sigma$	$\frac{1}{1+e^{-x}}$	$f(x)(1-f(x))$
Гиперболический тангенс $\tanh$	$\frac{e^x - e^{-x}}{e^x + e^{-x}}$	$1 - f^2(x)$
SoftSign	$\frac{x}{1+ x }$	$\frac{1}{(1+ x )^2}$
Ступенька (функция Хевисайда)	$\begin{cases} 0, & x < 0 \\ 1, & x \geq 0 \end{cases}$	0
SoftPlus	$\log(1 + e^x)$	$\frac{1}{1+e^{-x}}$
ReLU	$\begin{cases} 0, & x < 0 \\ x, & x \geq 0 \end{cases}$	$\begin{cases} 0, & x < 0 \\ 1, & x \geq 0 \end{cases}$
Leaky ReLU, Parameterized ReLU	$\begin{cases} ax, & x < 0 \\ x, & x \geq 0 \end{cases}$	$\begin{cases} a, & x < 0 \\ 1, & x \geq 0 \end{cases}$
ELU	$\begin{cases} \alpha(e^x - 1), & x < 0 \\ x, & x \geq 0 \end{cases}$	$\begin{cases} f(x) + \alpha, & x < 0 \\ 1, & x \geq 0 \end{cases}$

# Вопрос оптимизации

$$\min_w \frac{1}{m} \sum_{i=1}^m L(\hat{y}(x_i, w), y_i)$$

$L$  - функция потерь (optimisation objective)

$\hat{y}$  - выход модели

$w$  - веса модели

$x_i$  - входные данные

$y_i$  - истинные значения (targets)

$i = 1, ..m$

$m$  - размер тренировочной выборки

# Вопрос оптимизации

$$\min_w \frac{1}{m} \sum_{i=1}^m L(\hat{y}(x_i, w), y_i)$$

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

MSE = mean squared error

$n$  = number of data points

$Y_i$  = observed values

$\hat{Y}_i$  = predicted values

$L$  - функция потерь (optimisation objective)

$\hat{y}$  - выход модели

$w$  - веса модели

$x_i$  - входные данные

$y_i$  - истинные значения (targets)

$i = 1, ..m$

$m$  - размер тренировочной выборки

# Вопрос оптимизации

$$\min_w \frac{1}{m} \sum_{i=1}^m L(\hat{y}(x_i, w), y_i)$$

$L$  - функция потерь (optimisation objective)

$\hat{y}$  - выход модели

$w$  - веса модели

$x_i$  - входные данные

$y_i$  - истинные значения (targets)

$i = 1, ..m$

$m$  - размер тренировочной выборки

$$H_p(q) = -\frac{1}{N} \sum_{i=1}^N y_i \cdot \log(p(y_i)) + (1 - y_i) \cdot \log(1 - p(y_i))$$

Binary Cross-Entropy / Log Loss

## Градиентный спуск, а как вычислять

$$w^{(t+1)} = w^{(t)} - \eta \nabla L$$

$\eta$  - learning rate (step size, шаг обучения)

$\nabla_w L$  - градиент  $L$  относительно  $w$

## Градиентный спуск, а как вычислять

$$w^{(t+1)} = w^{(t)} - \eta \nabla L \iff \boxed{w_1 = w_1 - \eta \frac{\partial f}{\partial w_1} \quad \dots \quad w_n = w_n - \eta \frac{\partial f}{\partial w_n}}$$

$\eta$  - learning rate (step size, шаг обучения)

$\nabla_w L$  - градиент  $L$  относительно  $w$

# Градиентный спуск, а как вычислять

$$y = w_0 + w_1x_1 + w_2x_2 + \dots + w_nx_n + \varepsilon$$

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

$$w_1 = w_1 - \eta \frac{\partial f}{\partial w_1} \quad \dots \quad w_n = w_n - \eta \frac{\partial f}{\partial w_n}$$



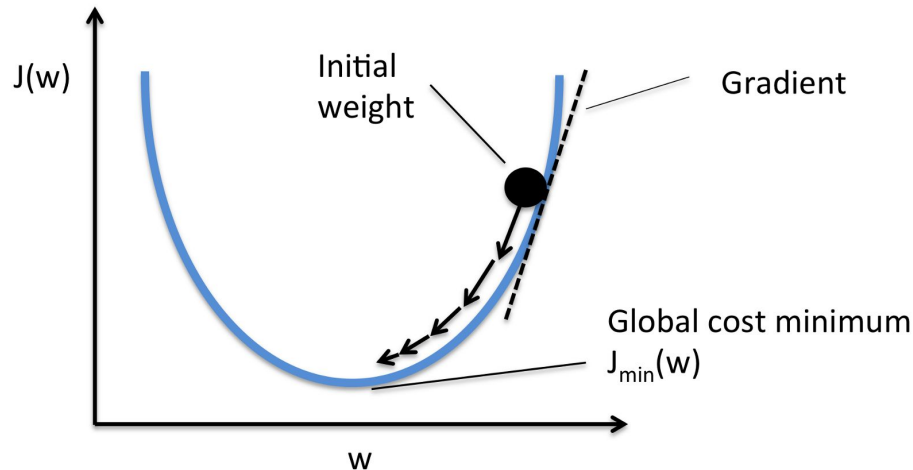


# Градиентный спуск, а как вычислять

$$y = w_1 x_1 + \varepsilon$$

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

$$w_1 = w_1 - \eta \frac{\partial f}{\partial w_1}$$



## Градиентный спуск, а как вычислять

$$y = wx + b$$

$$w = 1, \quad b = 0, \quad \eta = 0.01$$

$$(x = 3, y_{\text{true}} = 7)$$

$$\hat{y} = wx + b = 1 \cdot 3 + 0 = 3$$

$$L = (y_{\text{true}} - \hat{y})^2 = (7 - 3)^2 = 16$$

## Градиентный спуск, а как вычислять

$$\frac{\partial L}{\partial w} = 2(wx + b - y_{\text{true}})x$$

$$\frac{\partial L}{\partial b} = 2(wx + b - y_{\text{true}})$$

$$\frac{\partial L}{\partial w} = 2(3 - 7) \cdot 3 = -24$$

$$\frac{\partial L}{\partial b} = 2(3 - 7) = -8$$

## Градиентный спуск, а как вычислять

$$w = w - \eta \frac{\partial L}{\partial w} = 1 - 0.01 \cdot (-24) = 1 + 0.24 = 1.24$$

$$b = b - \eta \frac{\partial L}{\partial b} = 0 - 0.01 \cdot (-8) = 0 + 0.08 = 0.08$$

$$w = 1.24, \quad b = 0.08$$

## Градиентный спуск, а как вычислять

$$y = wx + b$$

$$w = 1.24, \quad b = 0.08, \quad \eta = 0.01$$

$$(x = 3, y_{\text{true}} = 7)$$

$$\hat{y} = wx + b = 1.24 \cdot 3 + 0.08 = 3.8$$

$$L = (y_{\text{true}} - \hat{y})^2 = (7 - 3.8)^2 = 10.24$$

# Теорема Цыбенко

Пусть  $\varphi$  любая непрерывная **сигмоидная функция**, например,  $\varphi(\xi) = 1/(1 + e^{-\xi})$ . Тогда, если дана любая непрерывная функция действительных переменных  $f$  на  $[0, 1]^n$  (или любое другое компактное подмножество  $\mathbb{R}^n$ ) и  $\varepsilon > 0$ , то существуют векторы  $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_N, \alpha$  и  $\theta$  и параметризованная функция  $G(\cdot, \mathbf{w}, \alpha, \theta) : [0, 1]^n \rightarrow \mathbb{R}$  такая, что для всех  $\mathbf{x} \in [0, 1]^n$  выполняется

$$|G(\mathbf{x}, \mathbf{w}, \alpha, \theta) - f(\mathbf{x})| < \varepsilon,$$

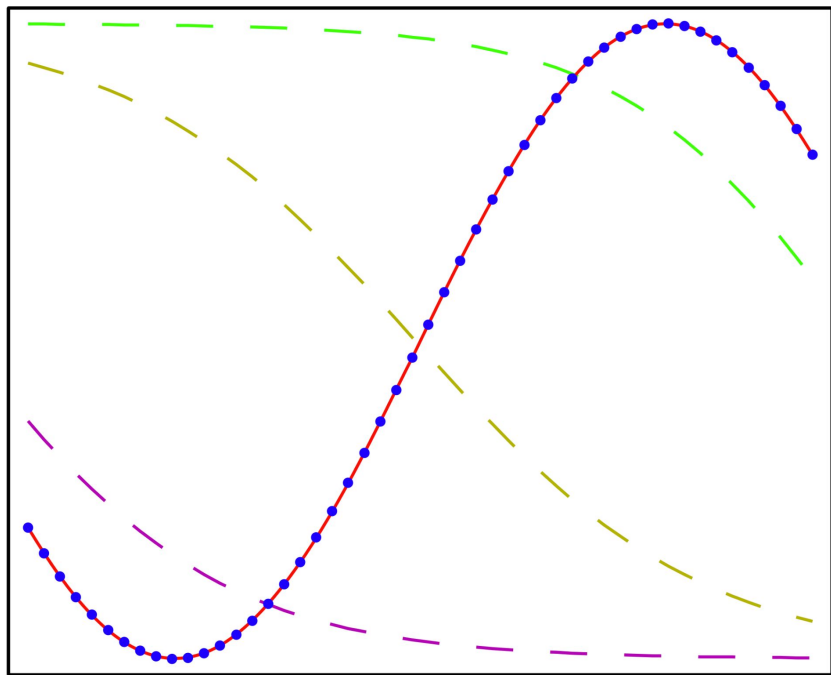
где

$$G(\mathbf{x}, \mathbf{w}, \alpha, \theta) = \sum_{i=1}^N \alpha_i \varphi(\mathbf{w}_i^T \mathbf{x} + \theta_i),$$

и  $\mathbf{w}_i \in \mathbb{R}^n$ ,  $\alpha_i, \theta_i \in \mathbb{R}$ ,  $\mathbf{w} = (\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_N)$ ,  $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_N)$ , и  $\theta = (\theta_1, \theta_2, \dots, \theta_N)$ .

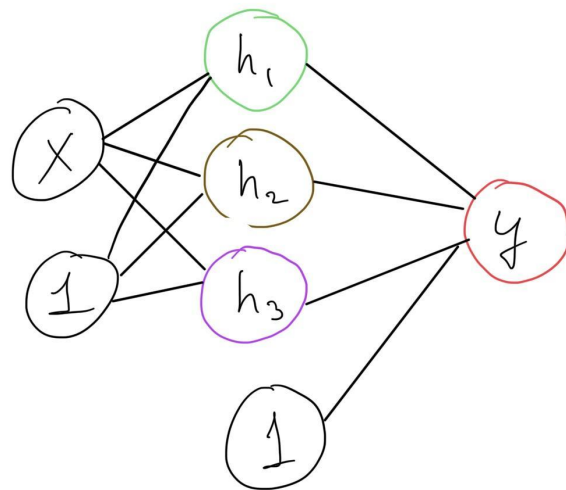
Cool, but impractical 😞

## More Toy Problems



$N = 50$  точек

MLP: 3 скрытых нейрона с  $\tanh$  активацией



## References

- Pattern recognition and machine learning (Bishop, 2006)
- Deep Learning (Goodfellow, 2014)
- Learning Theory from First Principles (Bach, 2023)
- Wikipedia
- Глубокое обучение (Николенко, 2018)
- Automatic Differentiation in Machine Learning: a Survey
- <https://github.com/MALINAYAGODA>
- <https://github.com/DaniilSergeev17>
- [github.com/karpathy/micrograd](https://github.com/karpathy/micrograd)

