

Skillbox

# Обучение нейронных сетей

# Градиентный спуск

# Нейронные сети

- Мы знаем, что такое нейронная сеть
- Знаем, как записать задачу оптимизации для нейронной сети

# Нейронные сети

- Мы знаем, что такое нейронная сеть
- Знаем, как записать задачу оптимизации для нейронной сети
- Но не знаем как ее решить )
  - Точнее, раньше мы подбирали коэффициенты руками, а это не очень практично

# Задача оптимизации

$$\mathcal{L} = \frac{1}{N} \sum_{i=1}^N L(y_i, f(\mathbf{x}_i, \boldsymbol{\theta})) \rightarrow \min_{\boldsymbol{\theta}}$$

$\mathbf{x}_i$  -- признаковое описание объекта  $i$

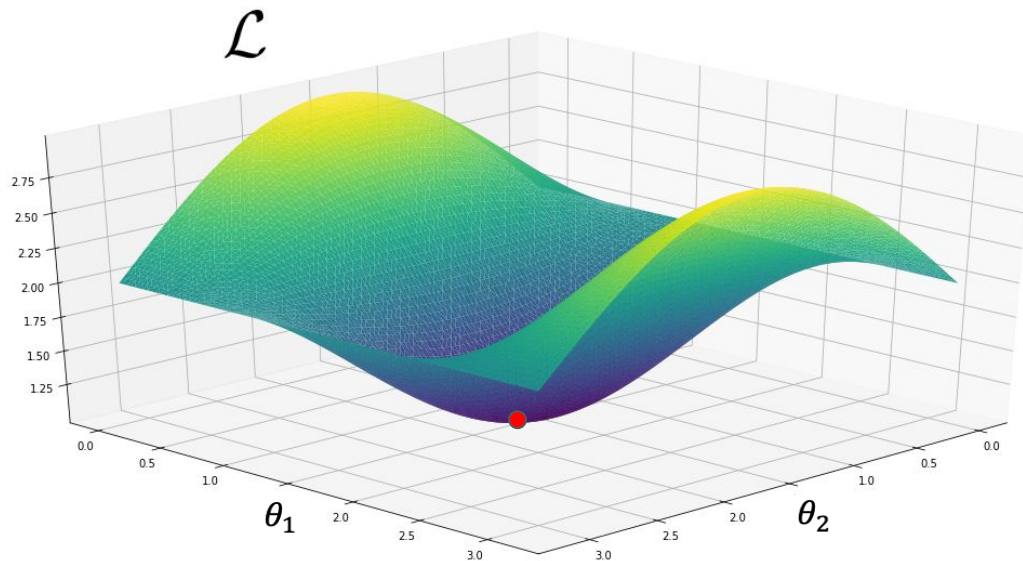
$y_i$  -- правильный ответ для объекта  $i$

$f$  -- наш алгоритм (нейронная сеть)

$\boldsymbol{\theta}$  -- параметры алгоритма

Нужно найти такие параметры, чтобы значение  $\mathcal{L}$  было минимальным

# Задача оптимизации

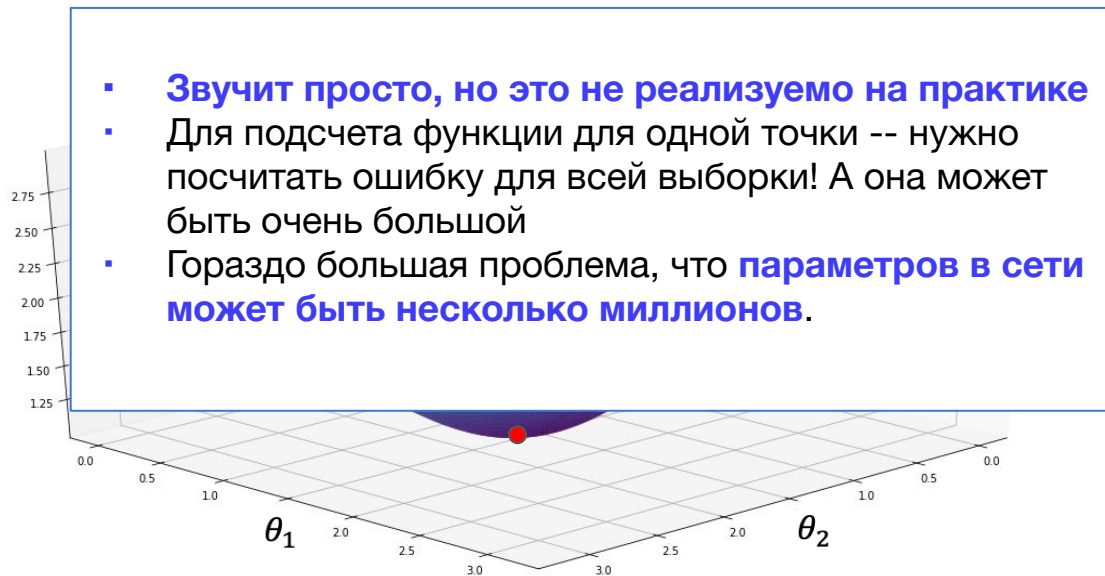


Два параметра (  $\theta_1$  ,  $\theta_2$  )

Посчитали значение  $\mathcal{L}$  для всех комбинаций параметров

Тогда ответом для нашей задачи будет набор параметров, соответствующий красной точке

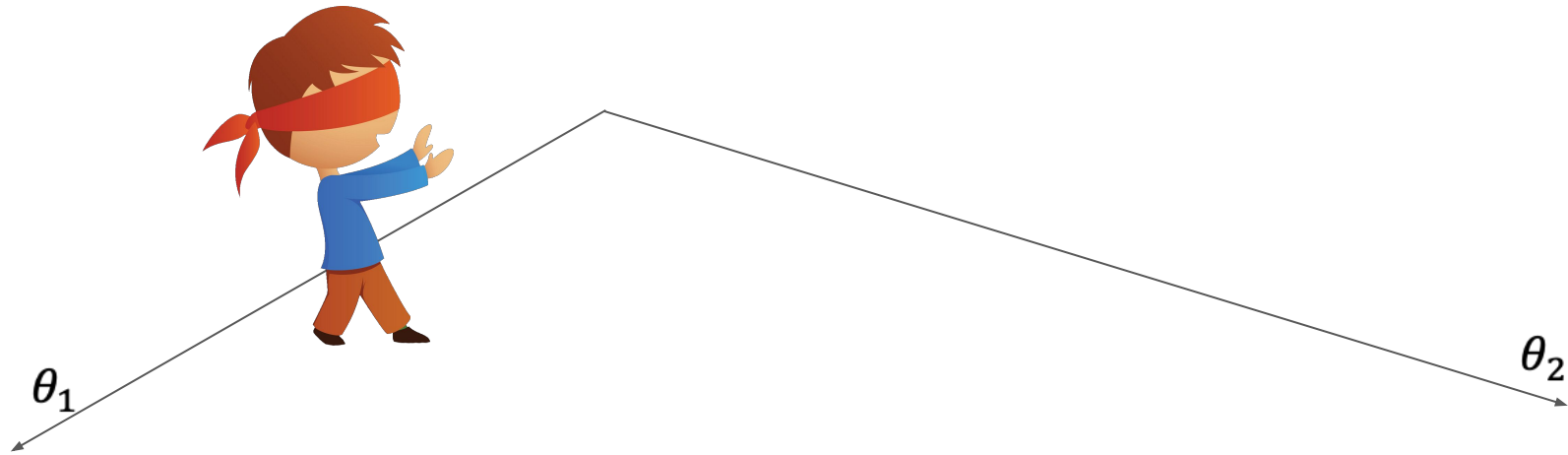
# Задача оптимизации



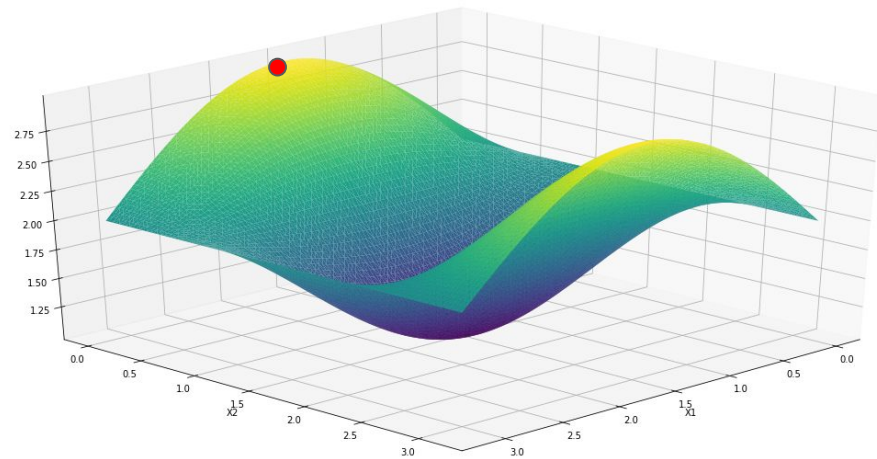
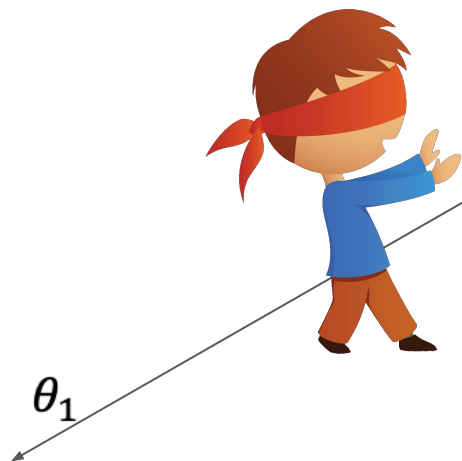
Два параметра (  $\theta_1$  ,  $\theta_2$  )

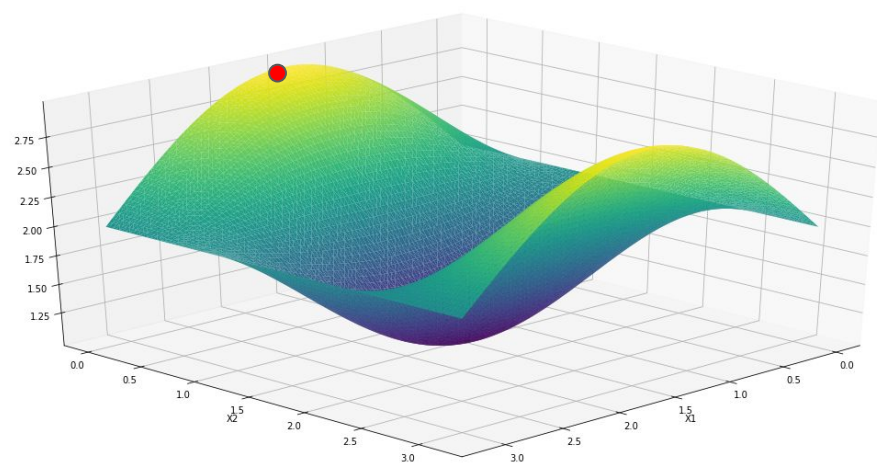
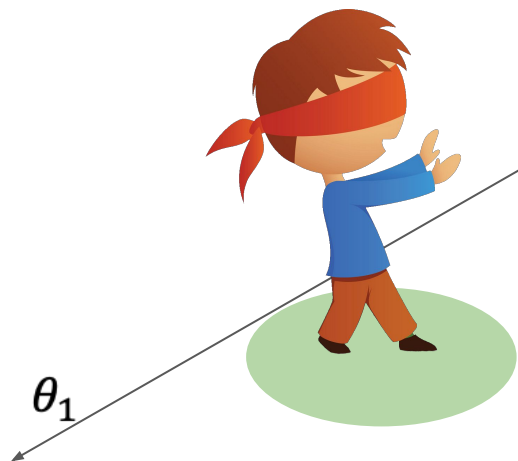
Посчитали значение  $\mathcal{L}$  для всех комбинаций параметров

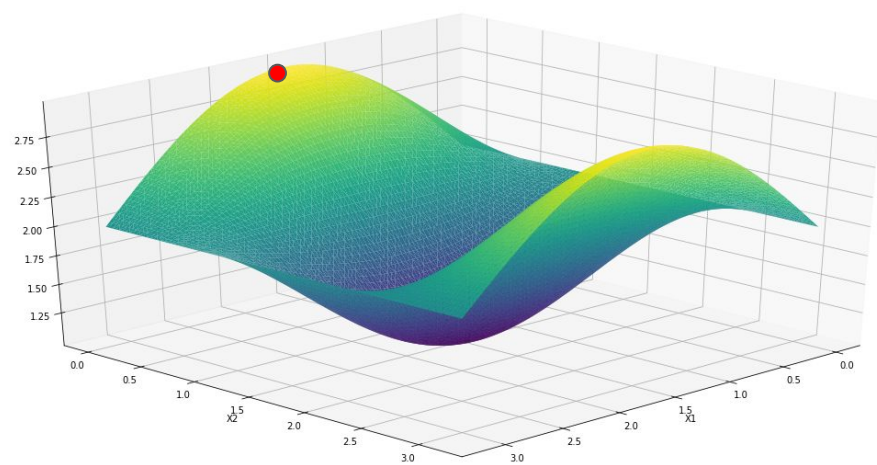
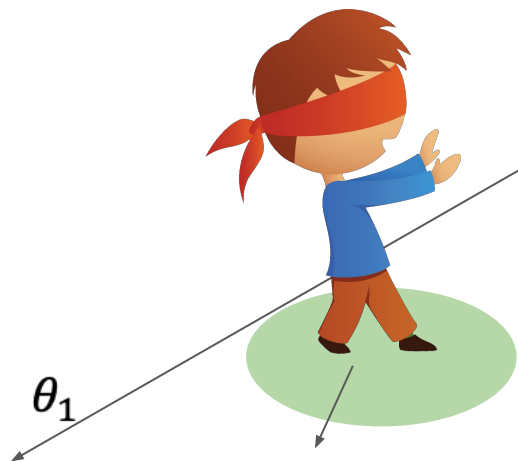
Тогда ответом для нашей задачи будет набор параметров, соответствующий красной точке



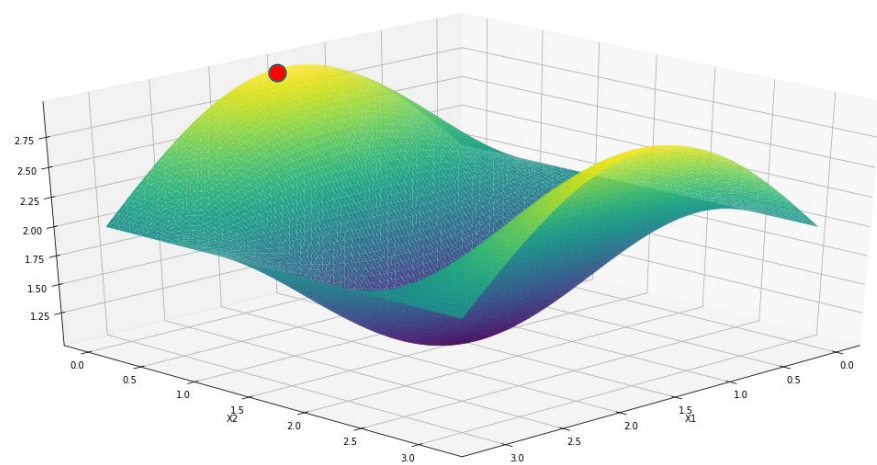
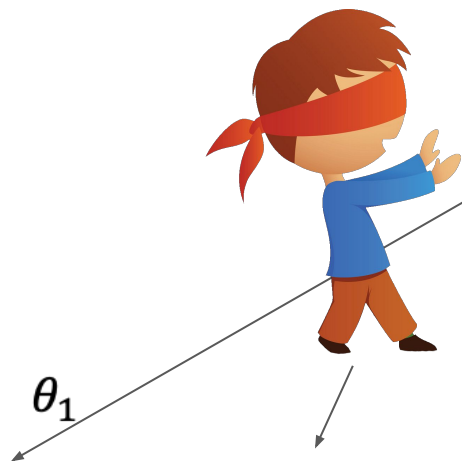


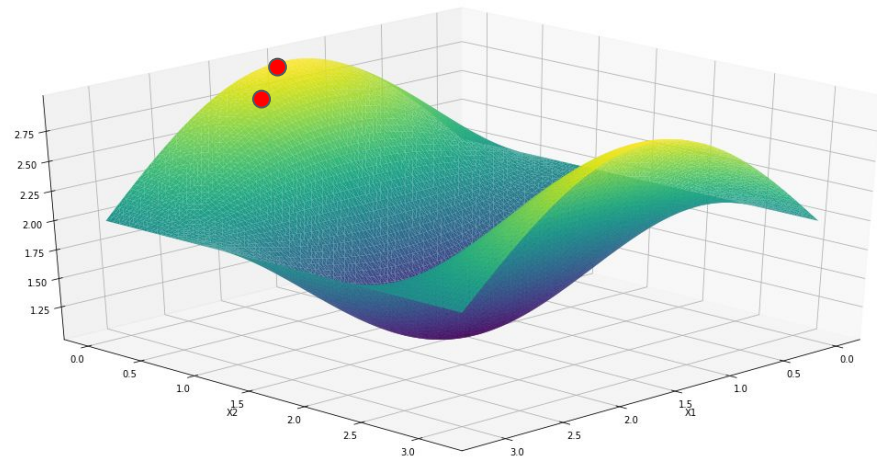
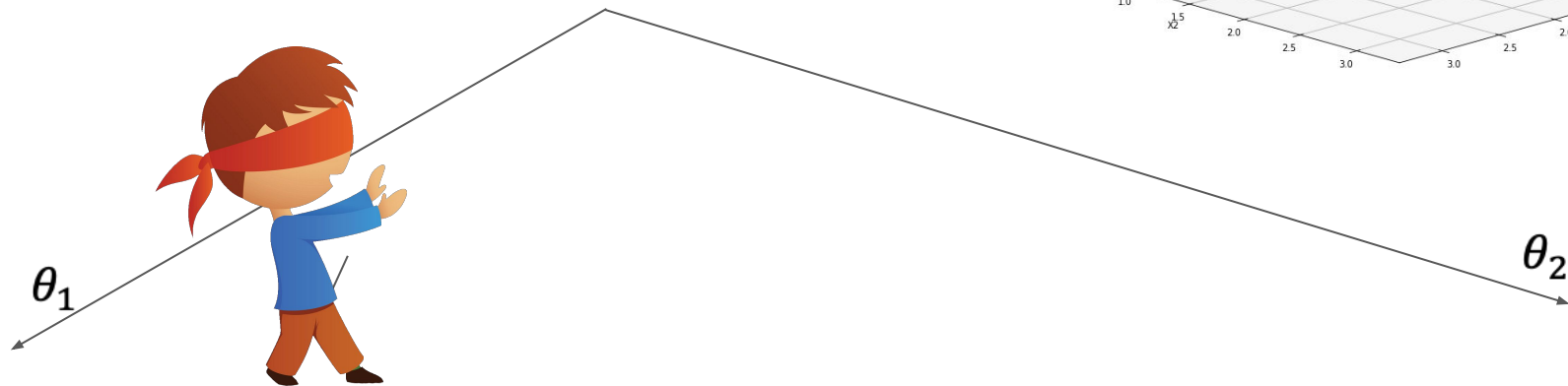


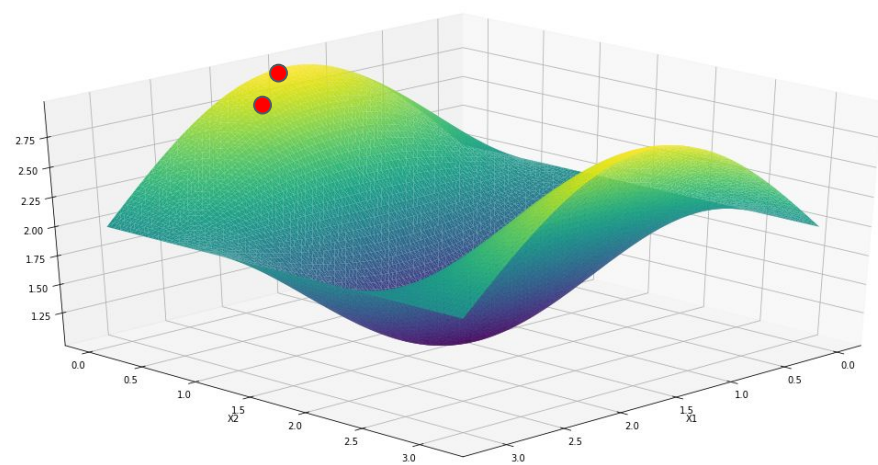
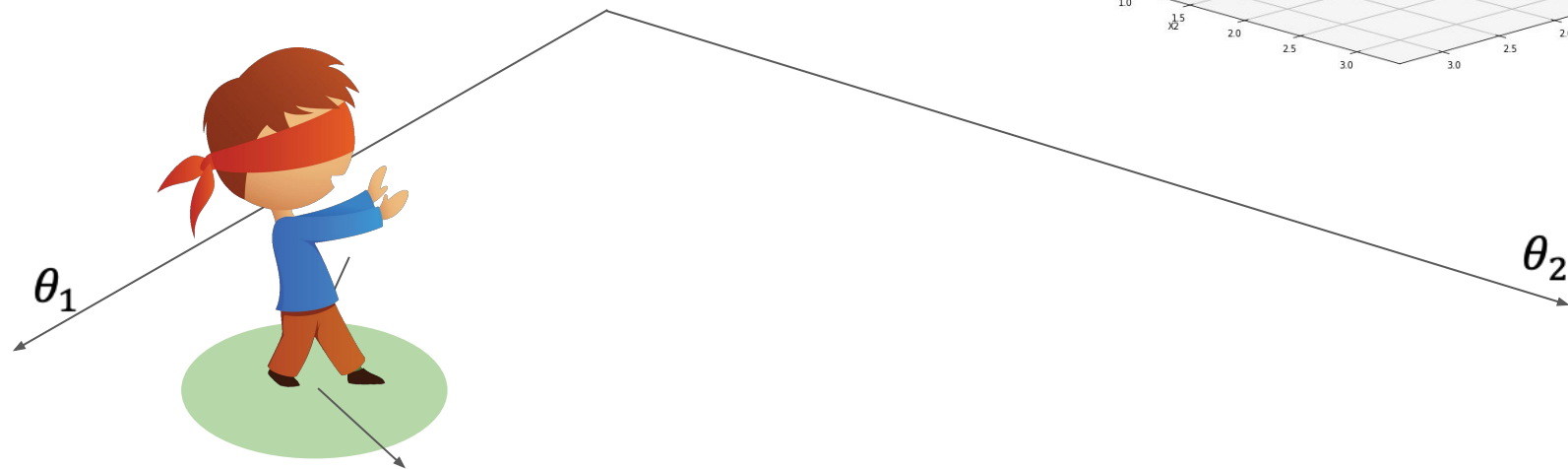


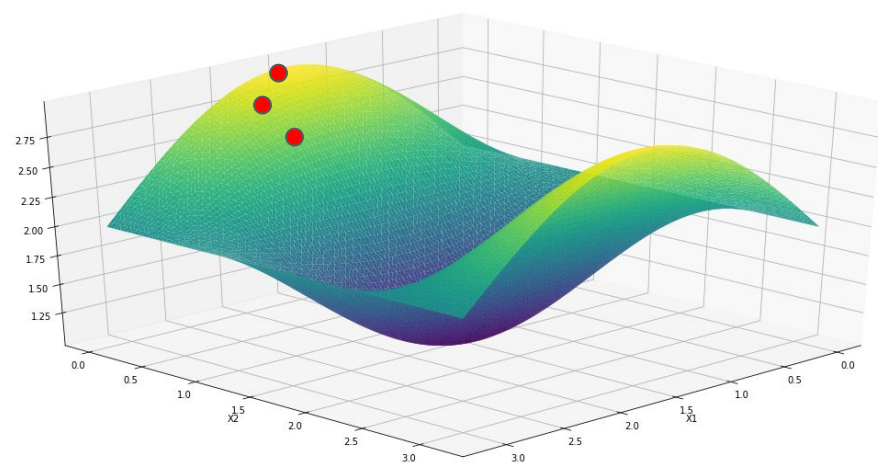
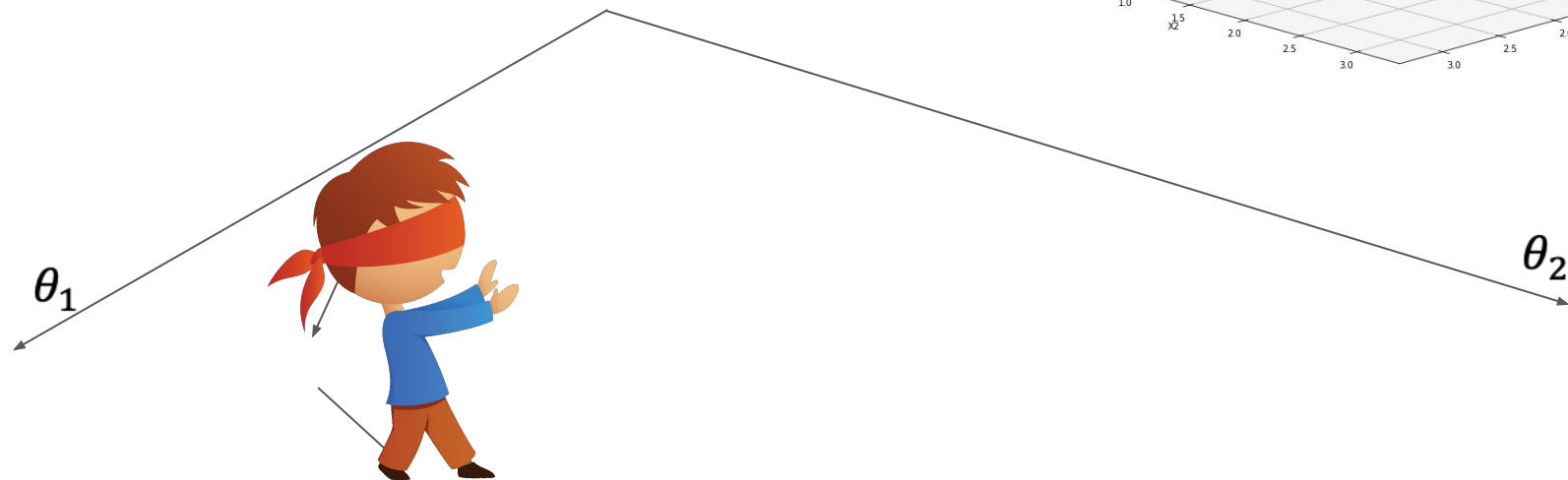


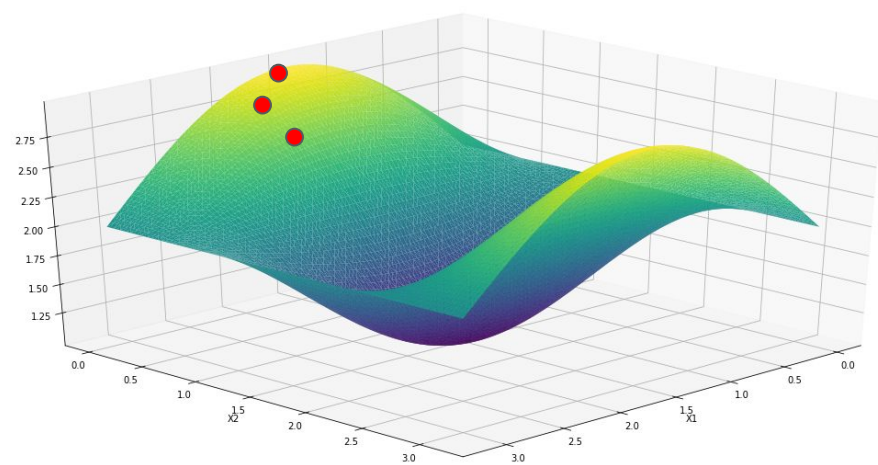
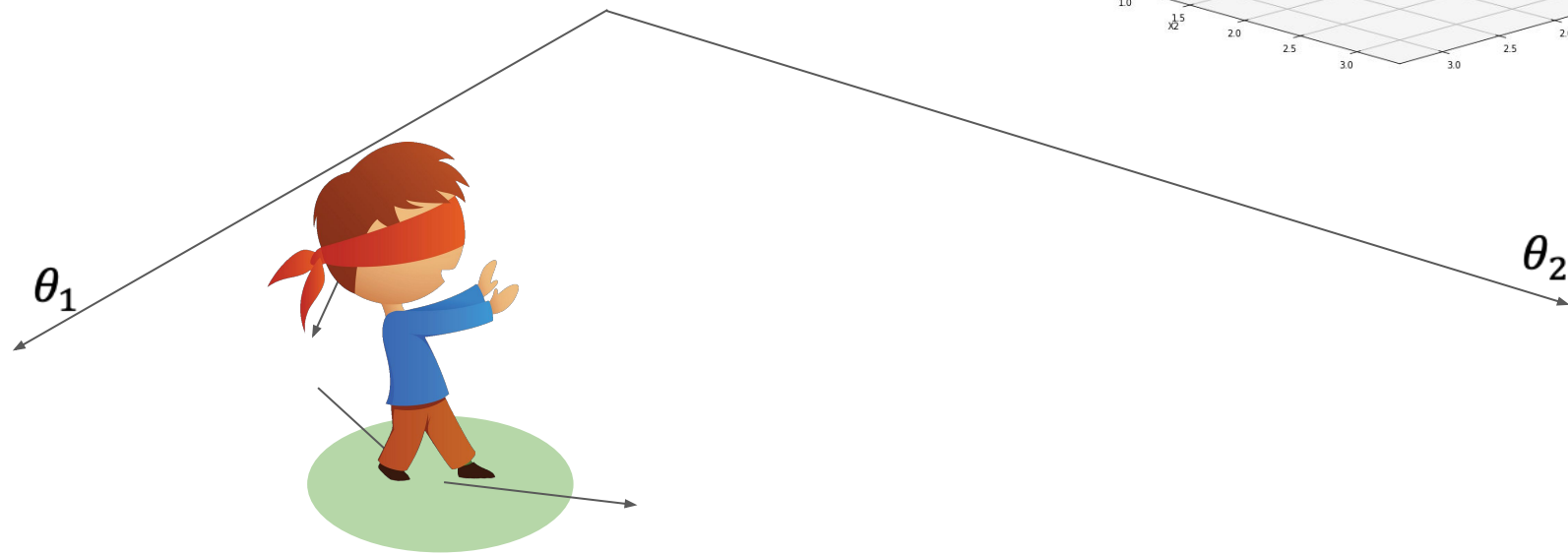
$\theta_2$



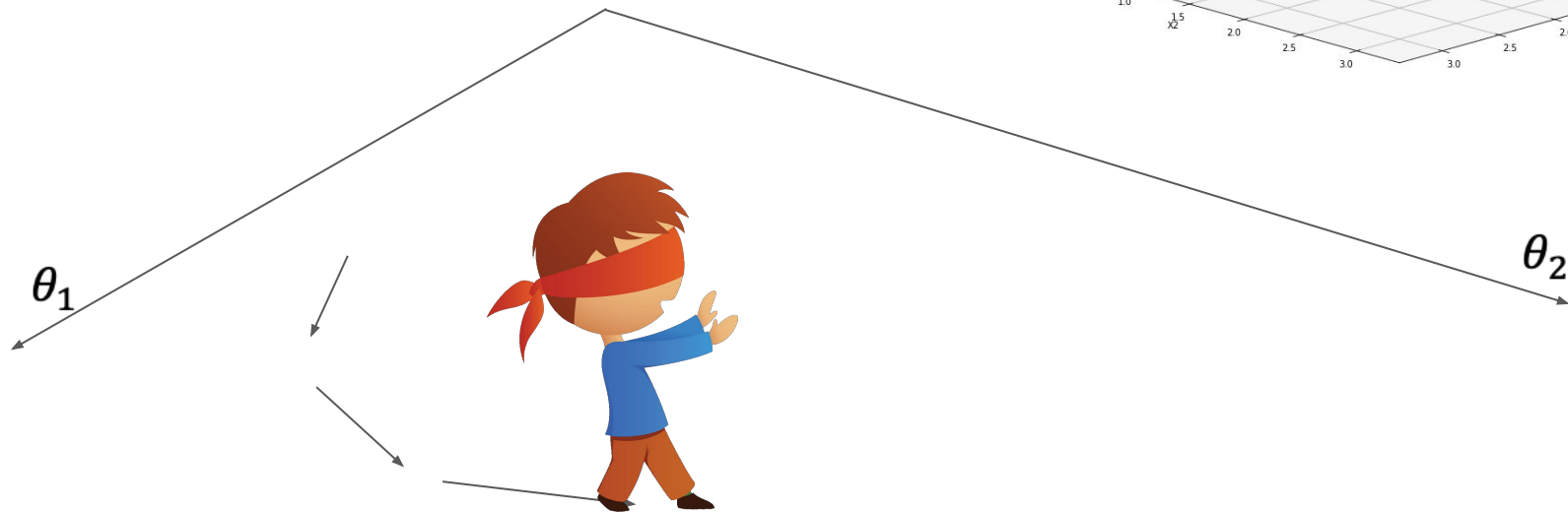
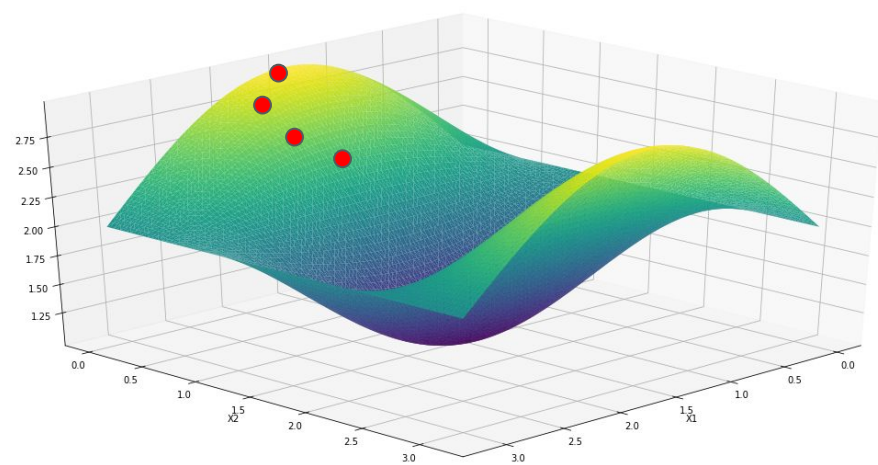


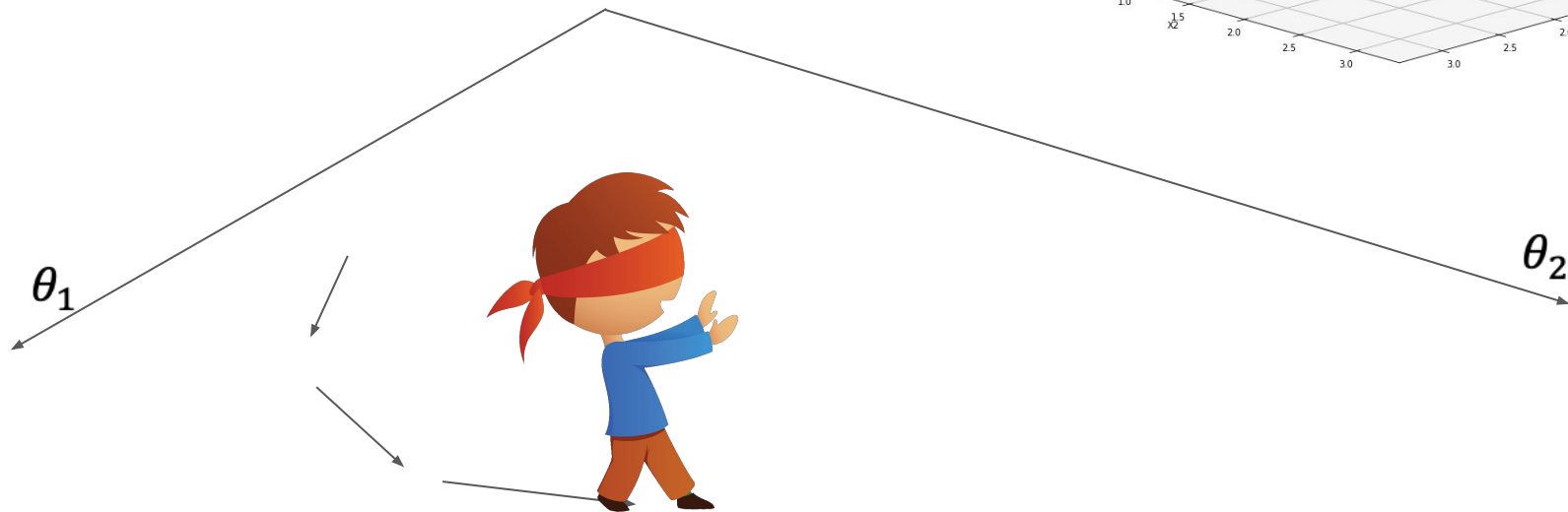
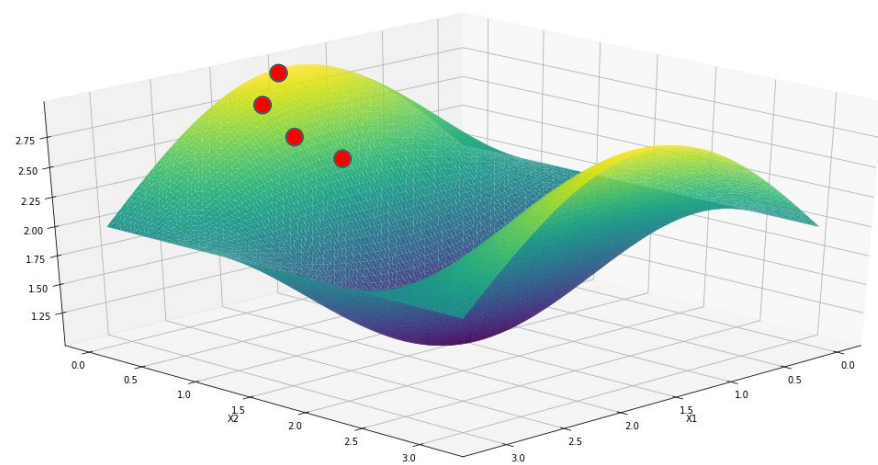


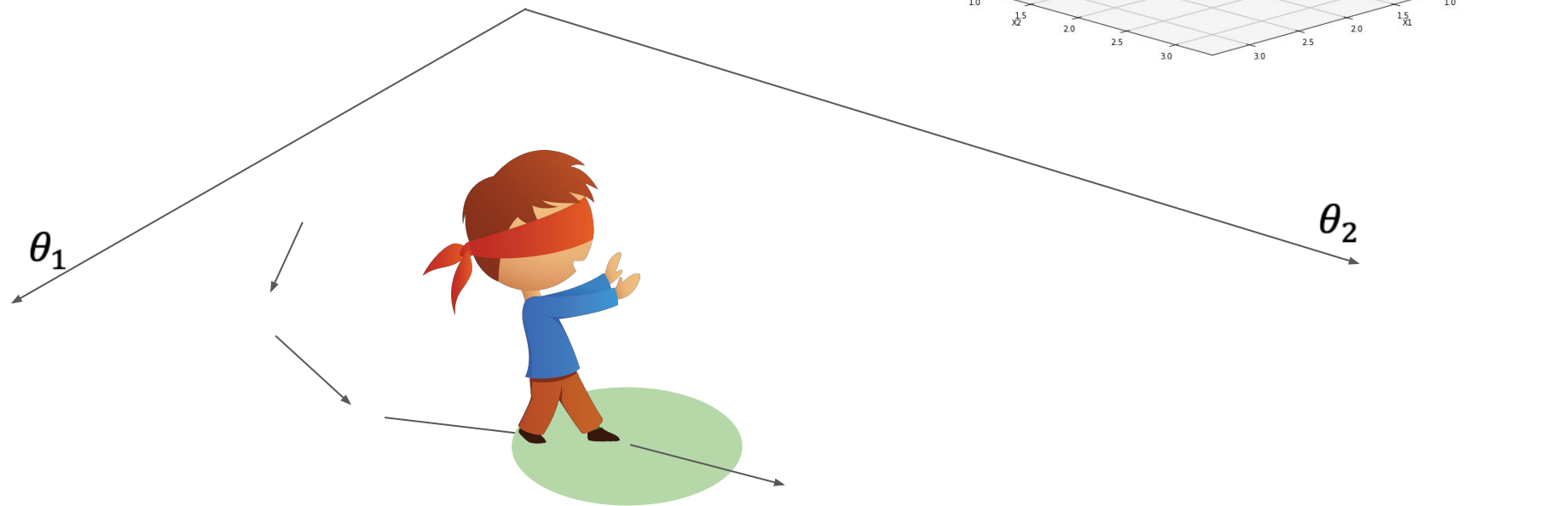


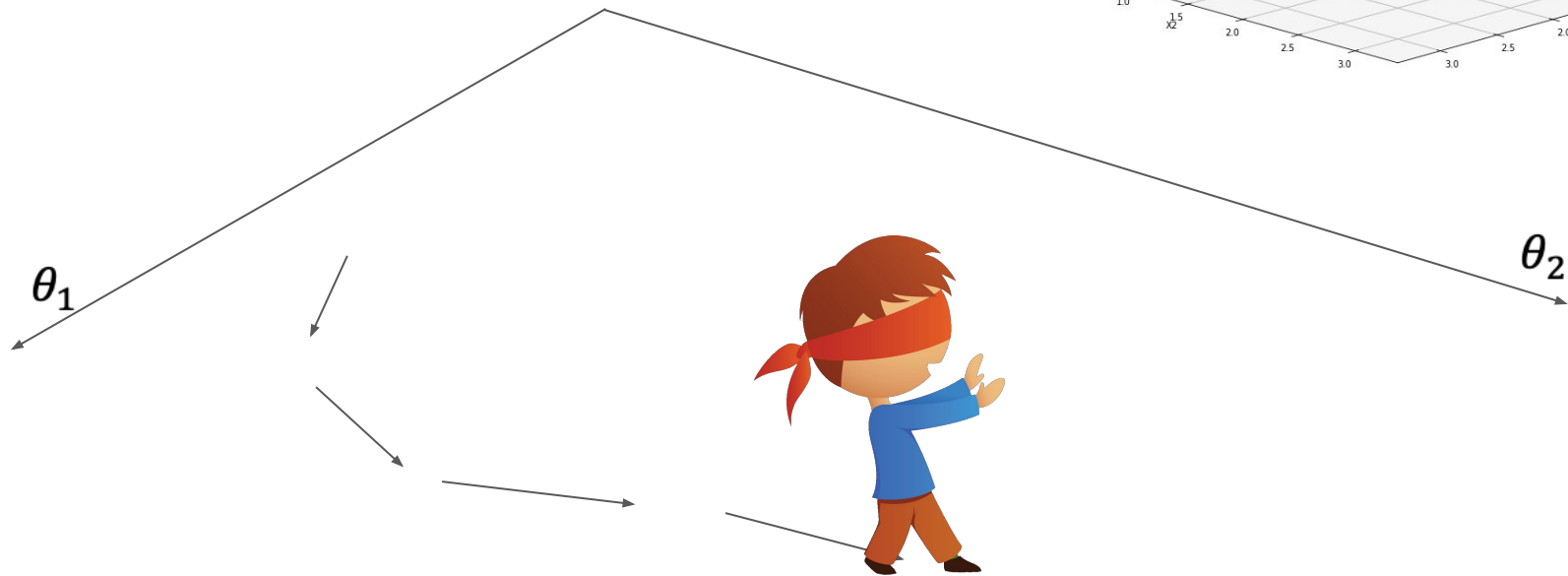
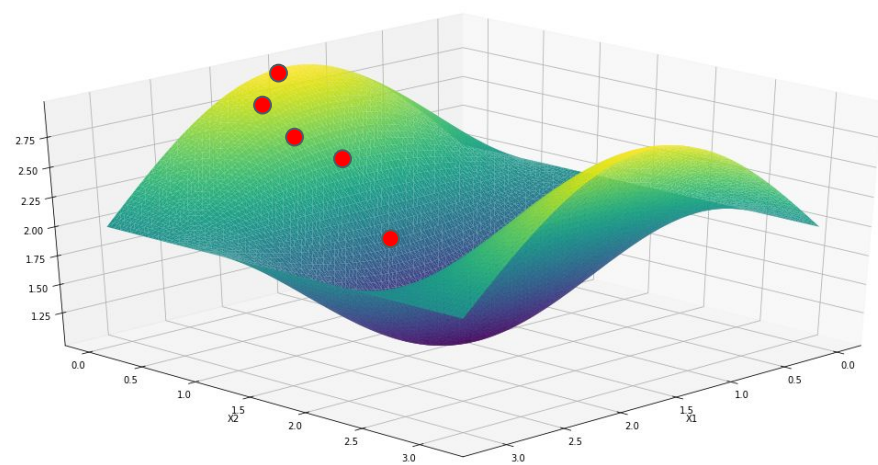


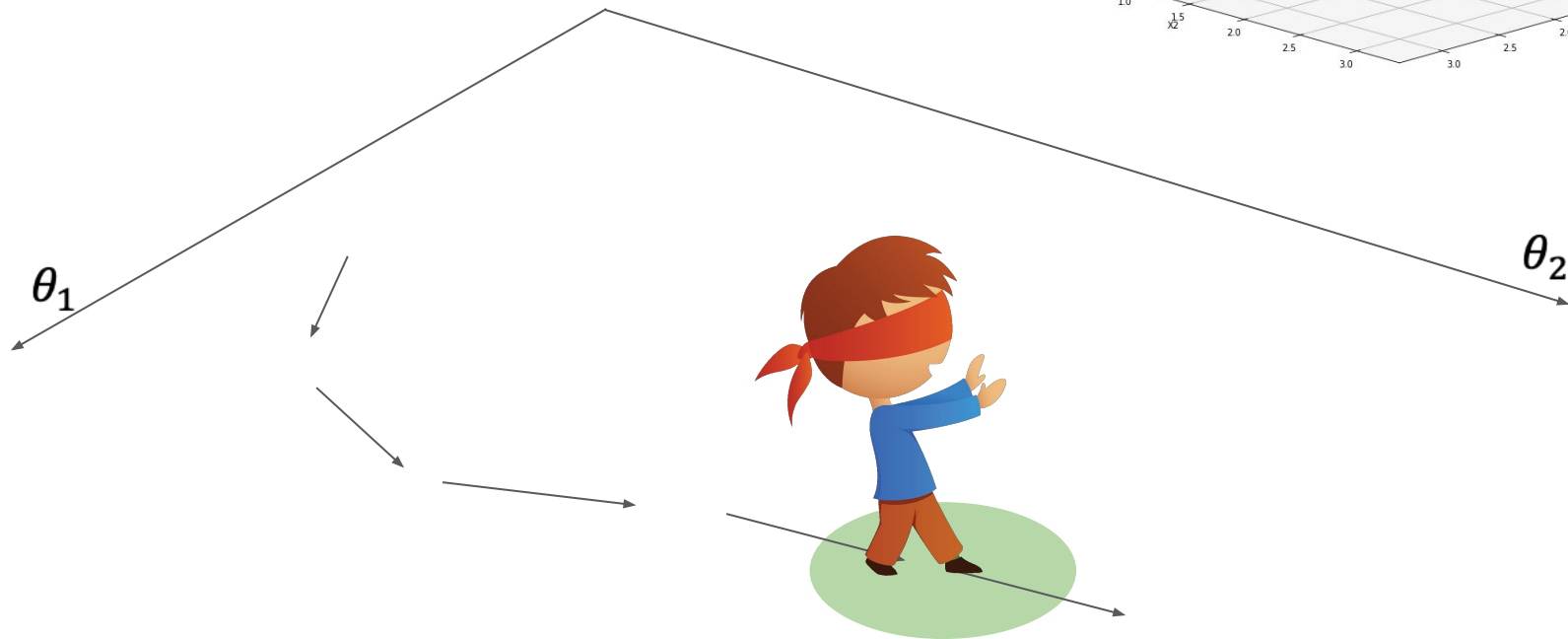
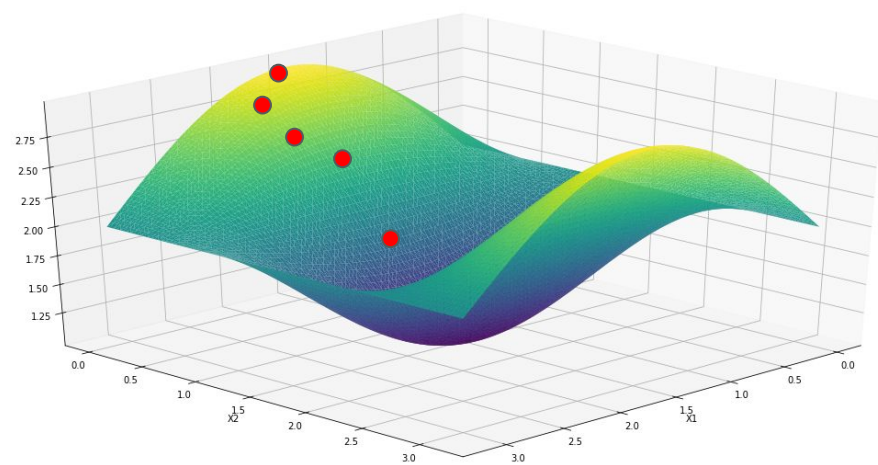


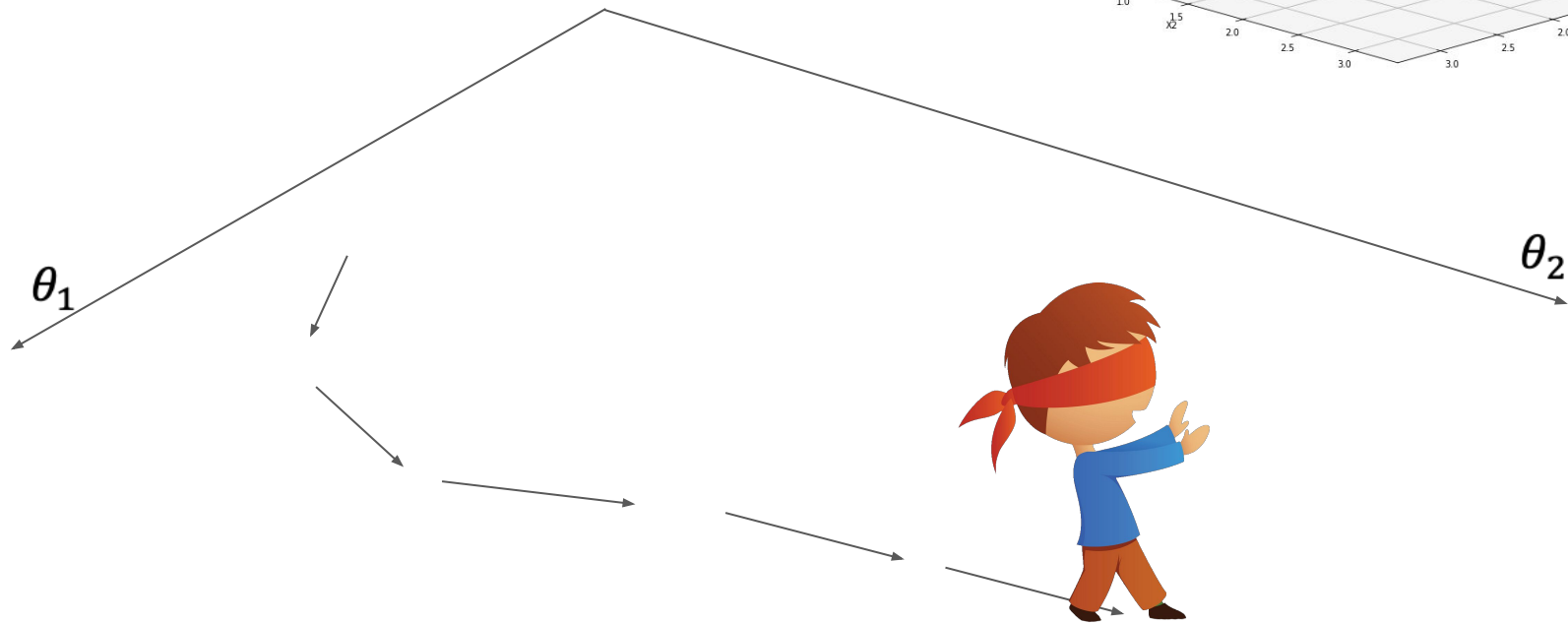
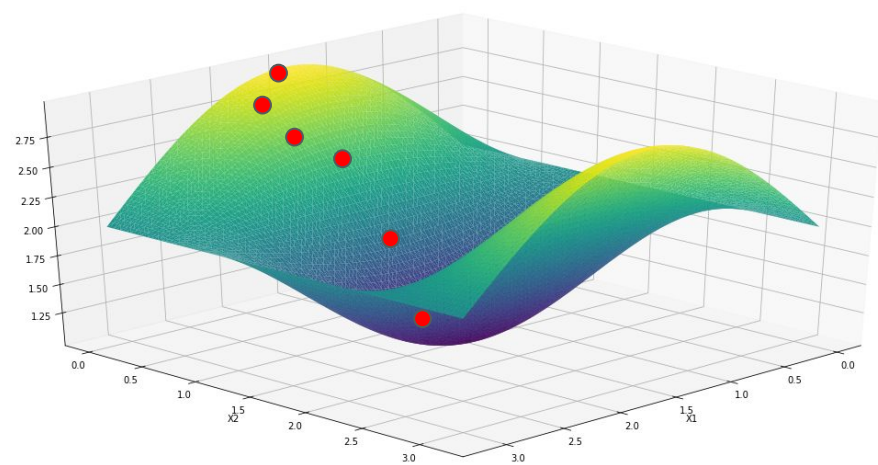


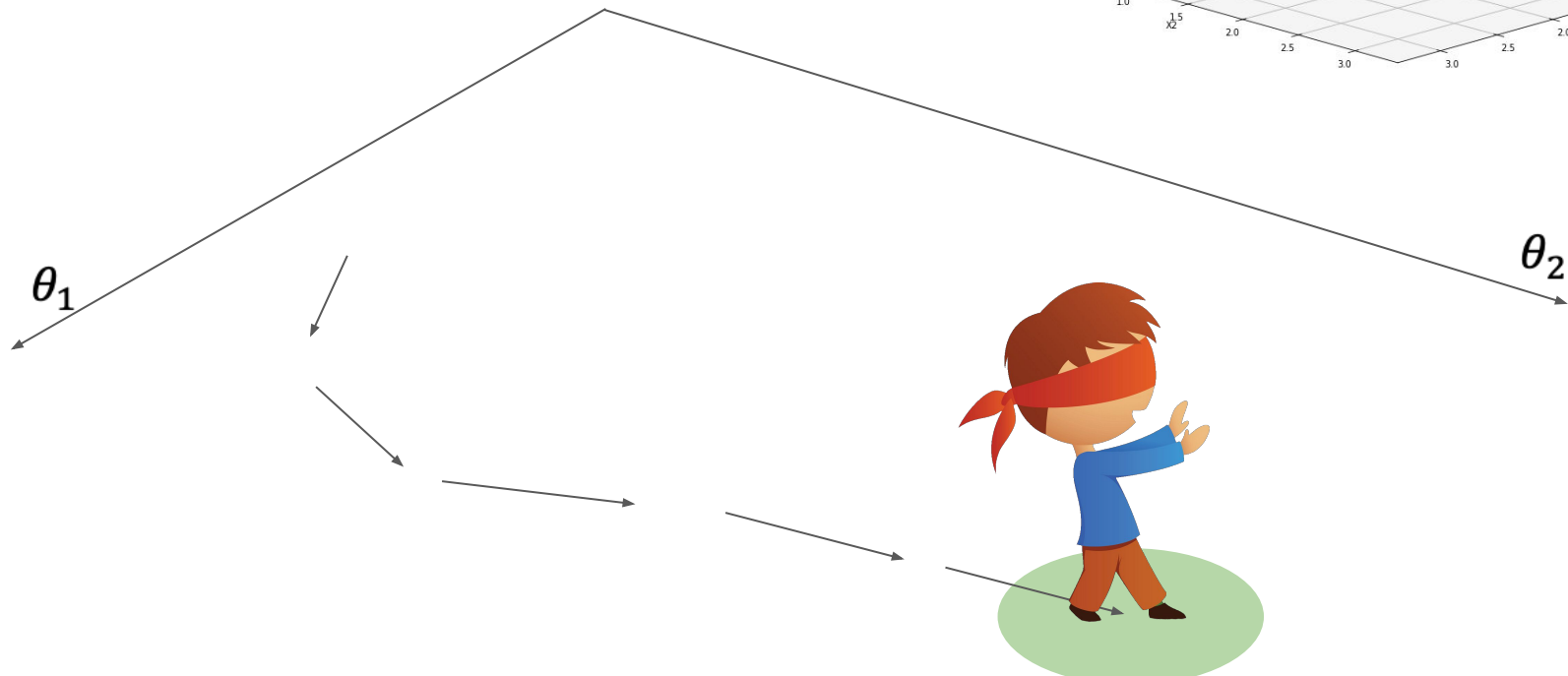
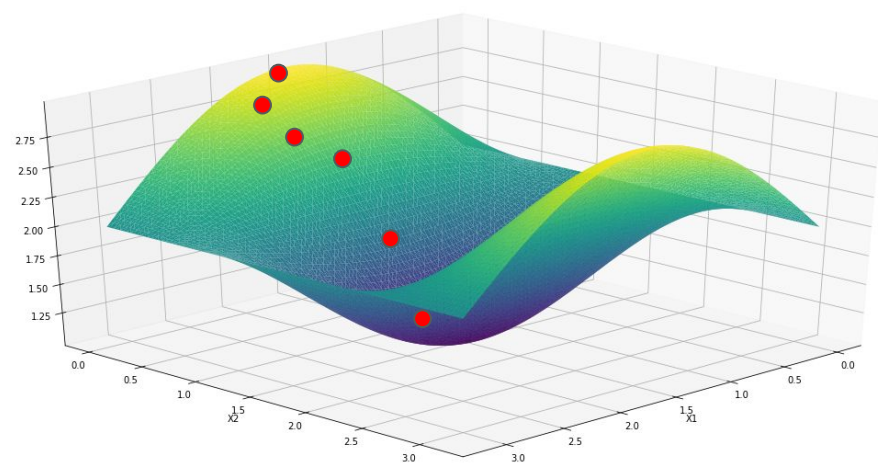


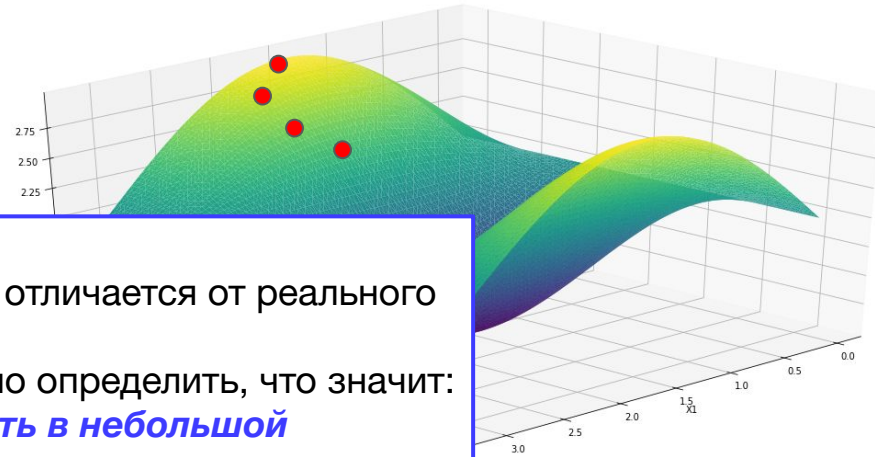




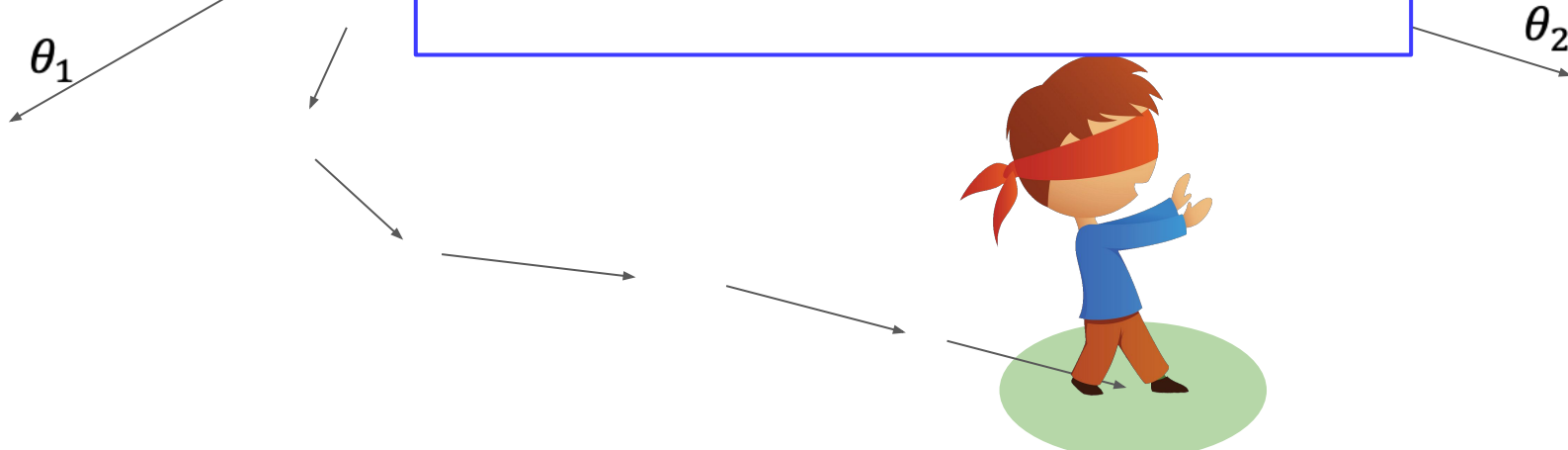








- Это описание не сильно отличается от реального алгоритма оптимизации!
- Нужно только формально определить, что значит: ***“прощупать поверхность в небольшой окрестности и выбрать направление в котором наблюдается наибольшее уменьшение высоты”***

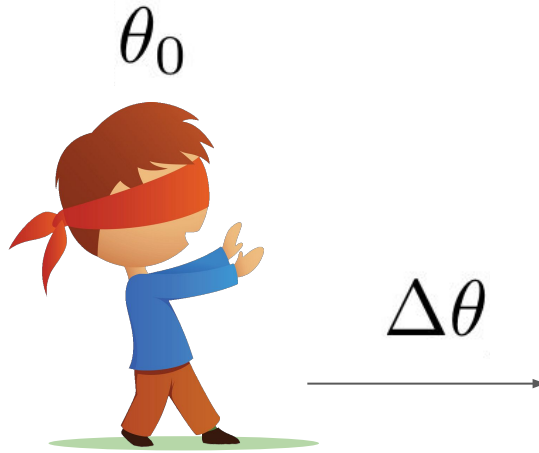


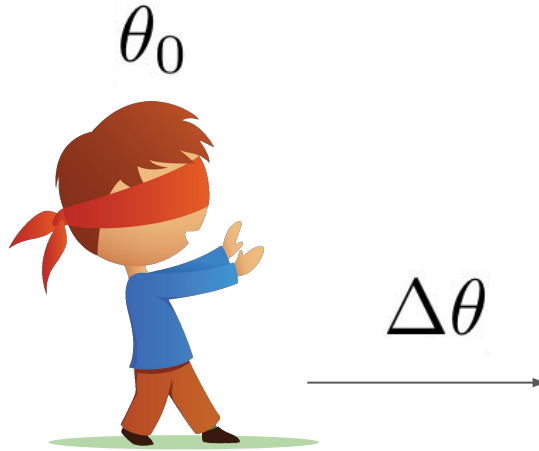










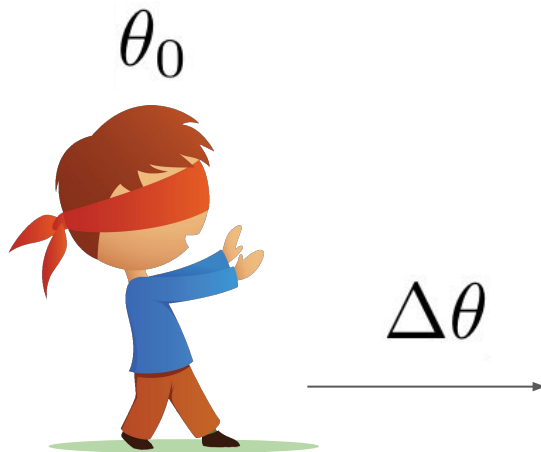


$$\frac{f(\theta_0 + \Delta\theta) - f(\theta_0)}{\Delta\theta}$$

$\Delta$ 

- $> 0$  : если пойти прямо, то значение функции увеличится
- $< 0$  : значение уменьшится
- **Нам нужно идти в сторону уменьшения**

$$\frac{f(\theta_0 + \Delta\theta) - f(\theta_0)}{\Delta\theta}$$

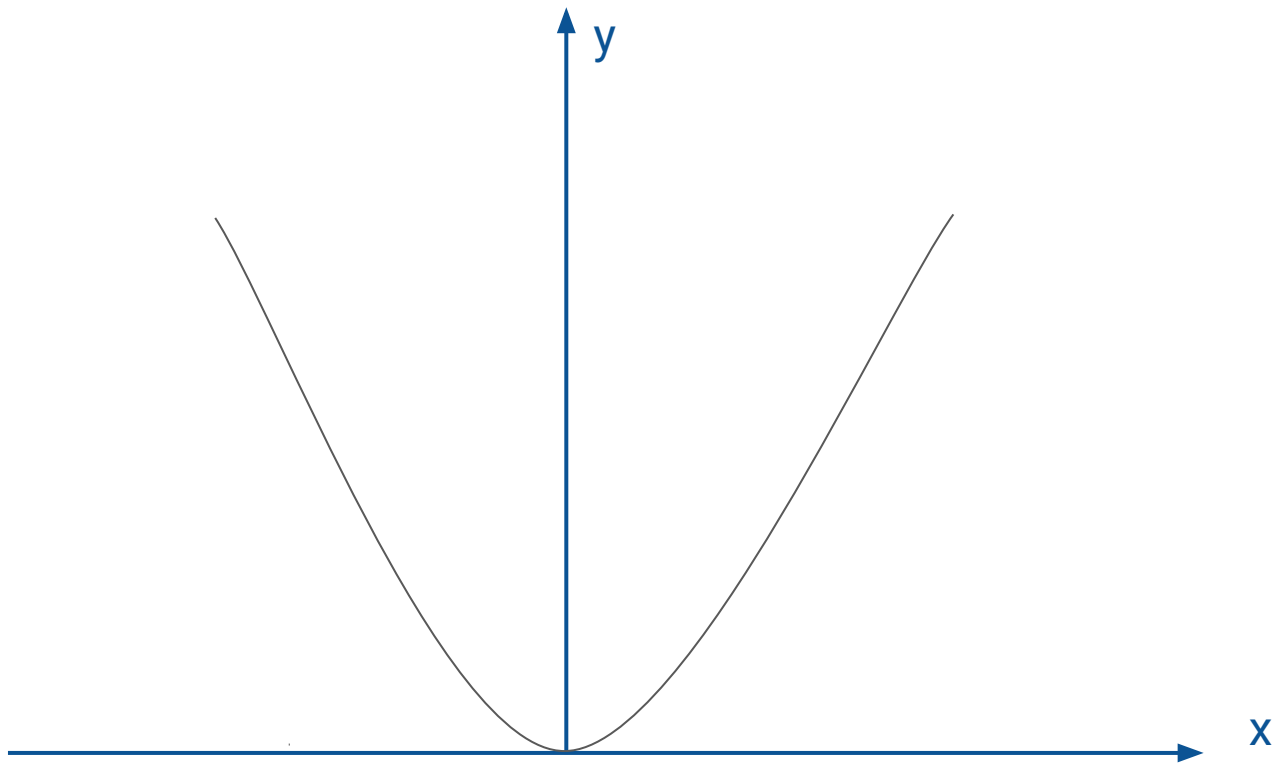


$$f'(\theta_0) = \lim_{\Delta\theta \rightarrow 0} \frac{f(\theta_0 + \Delta\theta) - f(\theta_0)}{\Delta\theta}$$

# Как найти минимум?

$$y = x^2$$

$$y' = 2x$$

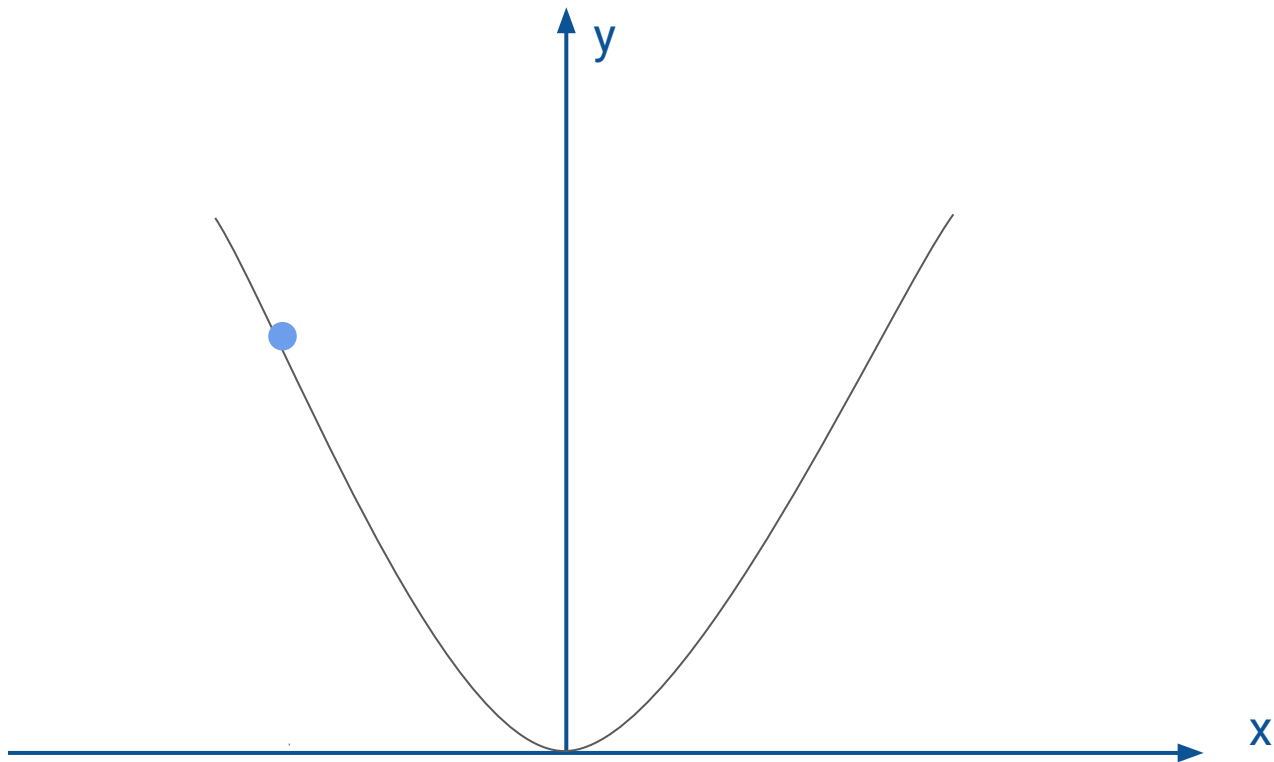




# Как найти минимум?

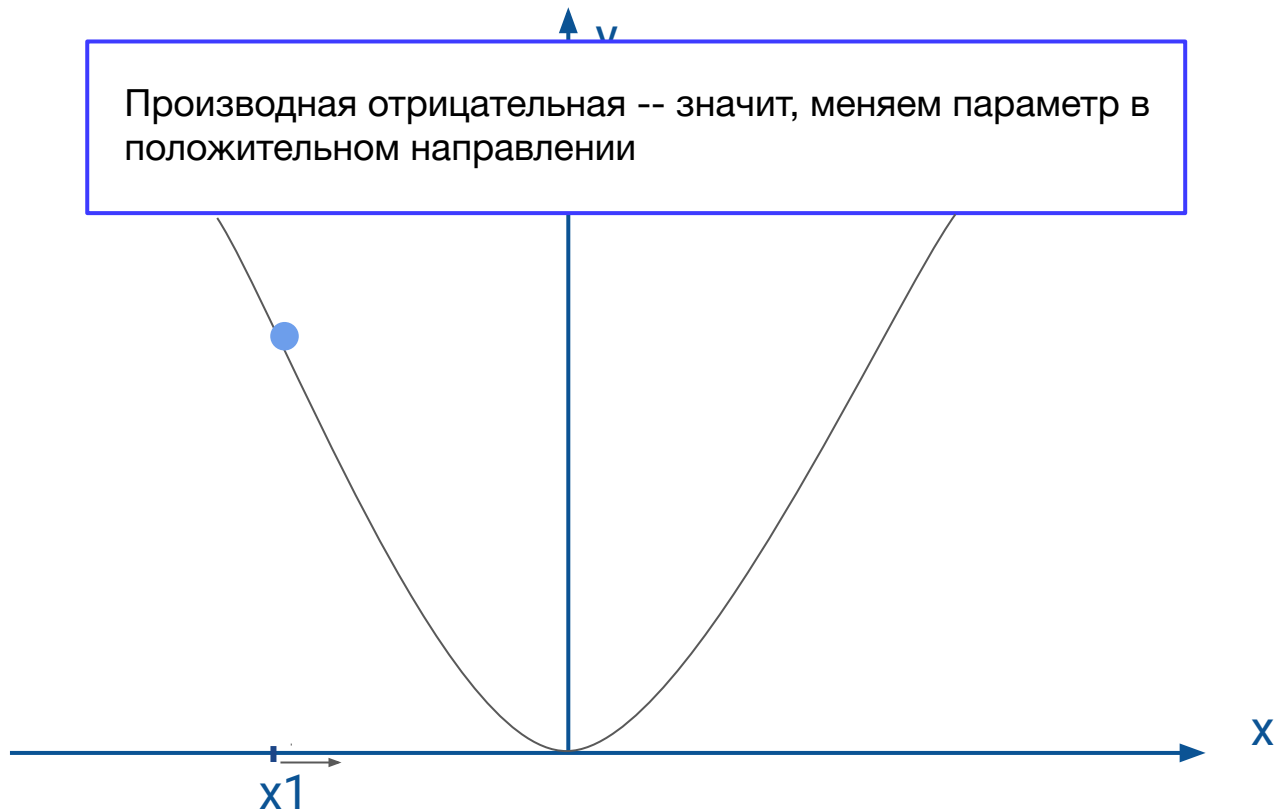
$$y = x^2$$

$$y' = 2x$$



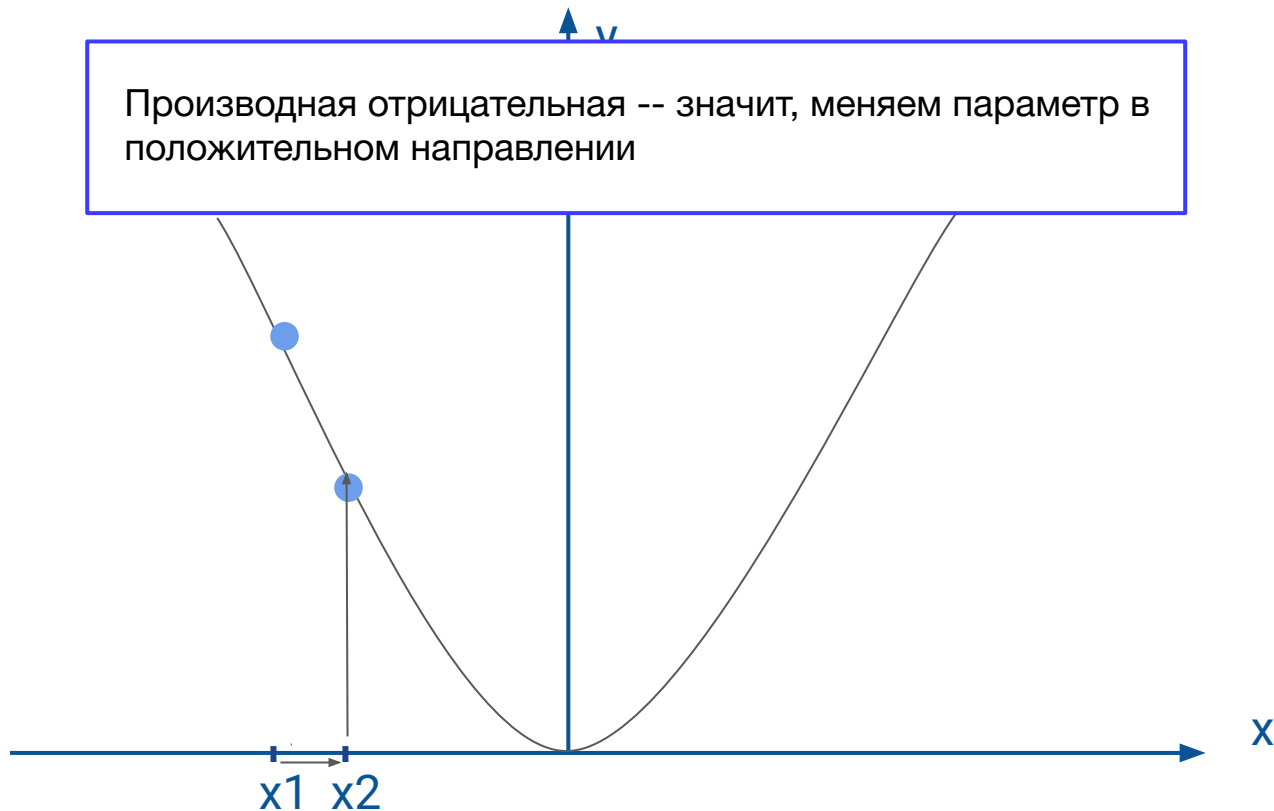
# Как найти минимум?

$$y = x^2$$
$$y' = 2x$$



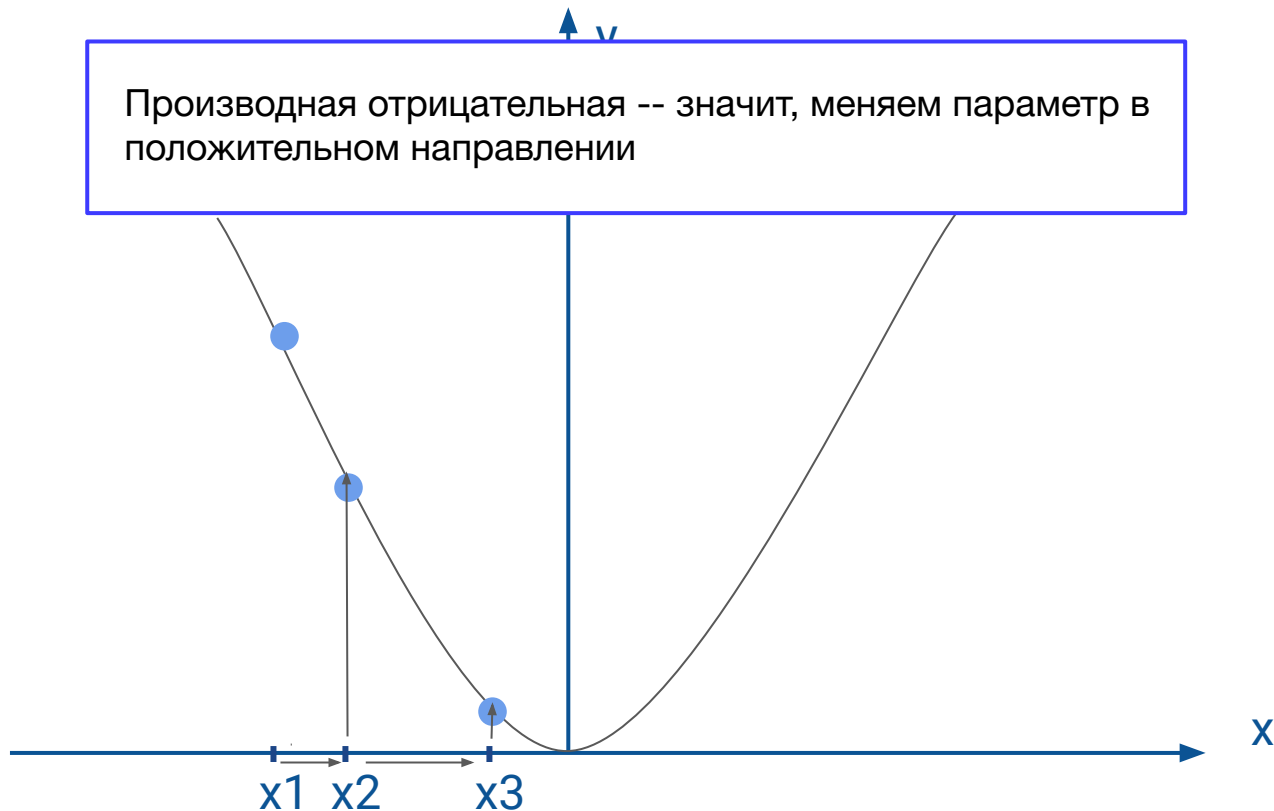
# Как найти минимум?

$$y = x^2$$
$$y' = 2x$$



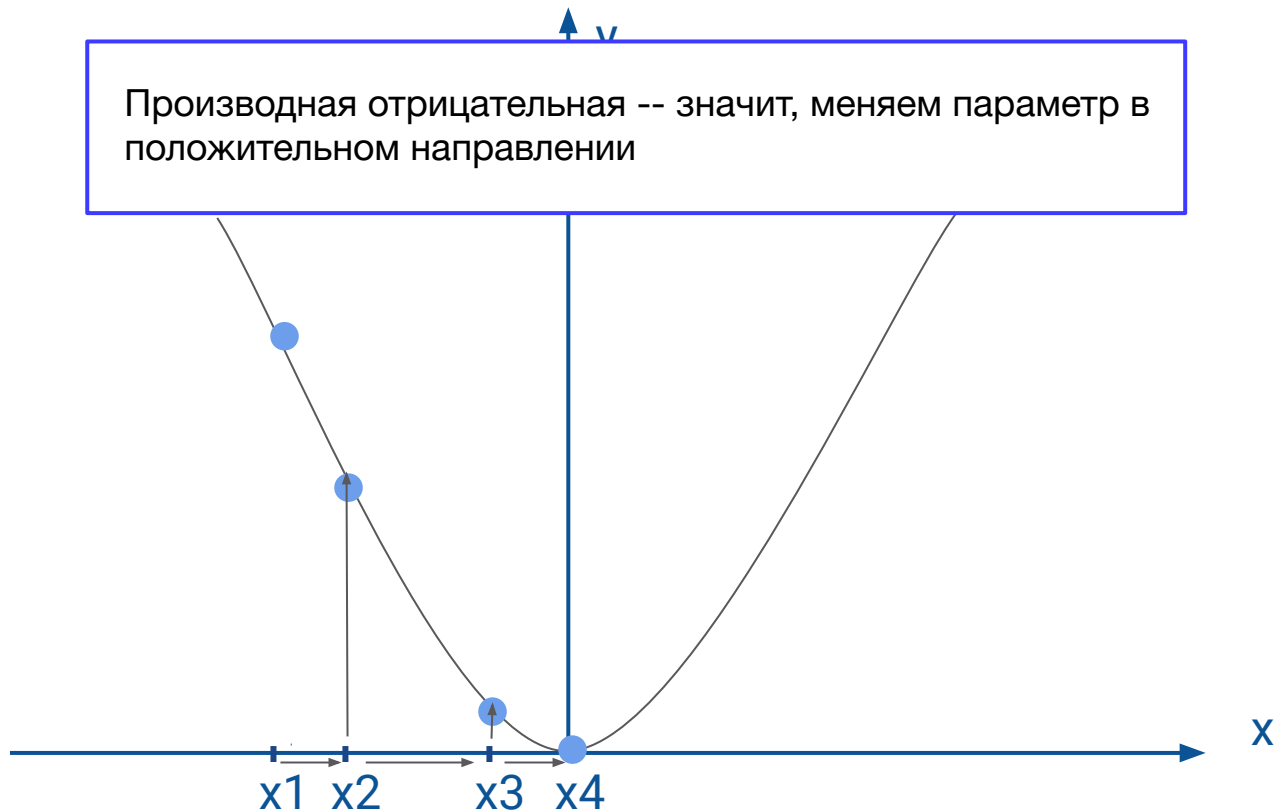
# Как найти минимум?

$$y = x^2$$
$$y' = 2x$$



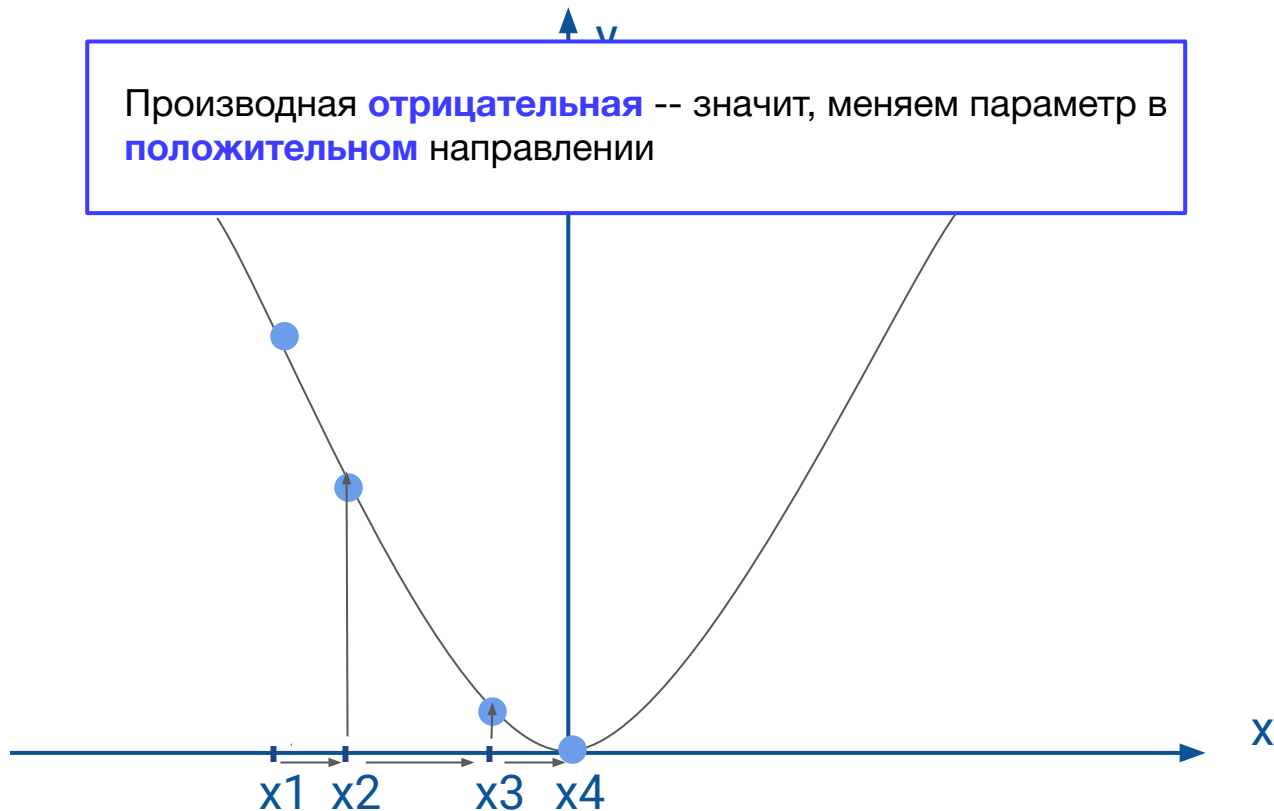
# Как найти минимум?

$$y = x^2$$
$$y' = 2x$$



# Как найти минимум?

$$y = x^2$$
$$y' = 2x$$



# Градиентный спуск

1. Инициализировать параметры случайным образом (выбрать  $\theta_0$ )
2. Вычислить значение функционала  $\mathcal{L}$
3. Вычислить  $\frac{\partial \mathcal{L}}{\partial \theta}(\theta_{i-1})$  при текущем значении параметра

4. Сделать шаг оптимизации:

$$\theta_i = \theta_{i-1} - \lambda \frac{\partial \mathcal{L}}{\partial \theta}(\theta_{i-1})$$

**а. После этого шага параметры изменились!**

5. Повторять 2-4 пока  $\mathcal{L}$  не перестанет меняться

# Градиентный спуск

1. Инициализировать параметры случайным образом (выбрать  $\theta_0$ )
2. Вычислить значение функционала  $\mathcal{L}$

3. Вычислить  $\frac{\partial \mathcal{L}}{\partial \theta}(\theta_{i-1})$  при текущем значении параметра

4. Сделать шаг оптимизации:

$$\theta_i = \theta_{i-1} - \boxed{\lambda} \frac{\partial \mathcal{L}}{\partial \theta}(\theta_{i-1})$$

Скорость обучения  
(learning rate)

**а. После этого шага параметры изменились!**

5. Повторять 2-4 пока  $\mathcal{L}$  не перестанет меняться



# Градиентный спуск

1. Инициализировать параметры случайным образом (выбрать  $\theta_0$ )
2. Вычислить значение функционала  $\mathcal{L}$
3. Вычислить  $\frac{\partial \mathcal{L}}{\partial \theta}(\theta_{i-1})$  при текущем значении параметра

4. Сделать шаг оптимизации:

$$\theta_i = \theta_{i-1} - \lambda \frac{\partial \mathcal{L}}{\partial \theta}(\theta_{i-1})$$

**а. После этого шага параметры изменились!**

5. Повторять 2-4 пока  $\mathcal{L}$  не перестанет меняться

# Градиентный спуск

1. Инициализировать параметры случайным образом (выбрать  $\theta_0$ )
2. Вычислить значение функционала  $\mathcal{L}$
3. Вычислить  $\frac{\partial \mathcal{L}}{\partial \theta}(\theta_{i-1})$  при текущем значении параметра
4. Сделать шаг оптимизации:

$$\theta_i = \theta_{i-1} - \lambda \frac{\partial \mathcal{L}}{\partial \theta}(\theta_{i-1})$$

**а. После этого шага параметры изменились!**

5. Повторять 2-4 пока  $\mathcal{L}$  не перестанет меняться

*\* градиент -- вектор из производных по каждому из параметров (для многомерного случая)*

# Градиентный спуск

1. Инициализировать параметры случайным образом (выбрать  $\theta_0$ )
2. Вычислить значение функционала  $\mathcal{L}$
3. Вычислить  $\frac{\partial \mathcal{L}}{\partial \theta}(\theta_{i-1})$  при текущем значении параметра
4. Сделать шаг оптимизации:
$$\theta_i = \theta_{i-1} - \lambda \frac{\partial \mathcal{L}}{\partial \theta}(\theta_{i-1})$$
  - а. После этого шага параметры изменились!
5. Повторять 2-4 пока  $\mathcal{L}$  не перестанет меняться

“Обучение”  
алгоритма



*\* градиент -- вектор из производных по каждому из параметров (для многомерного случая)*

# Градиентный спуск

1. Инициализировать параметры случайным образом (выбрать  $\theta_0$ )
2. Вычислить значение функционала  $\mathcal{L}$
3. Вычислить  $\frac{\partial \mathcal{L}}{\partial \theta}(\theta_{i-1})$  при текущем значении параметра
4. Сделать шаг оптимизации:

$$\theta_i = \theta_{i-1} - \lambda \frac{\partial \mathcal{L}}{\partial \theta}(\theta_{i-1})$$

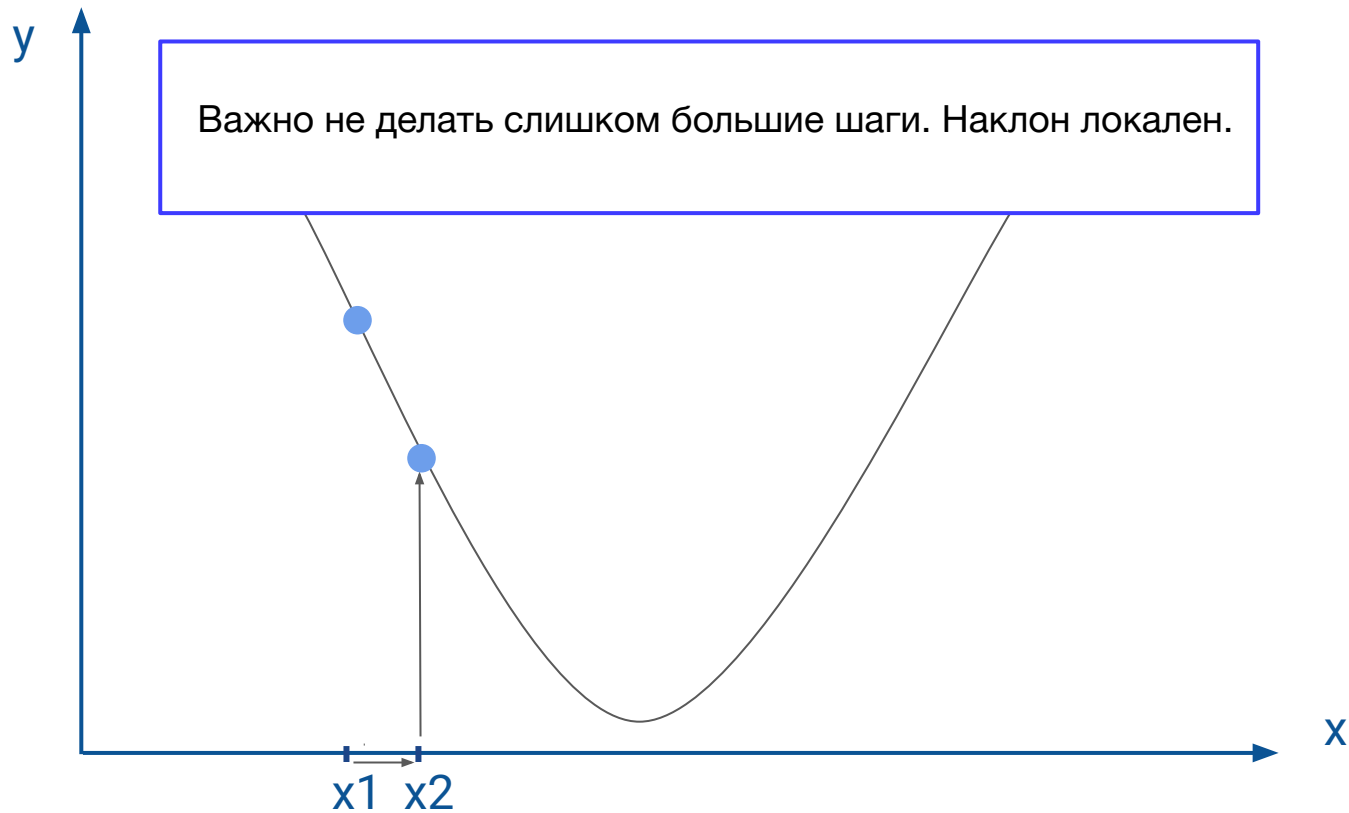
  - а. После этого шага параметры изменились!
5. Повторять 2-4 пока  $\mathcal{L}$  не перестанет меняться

“Обучение”  
алгоритма

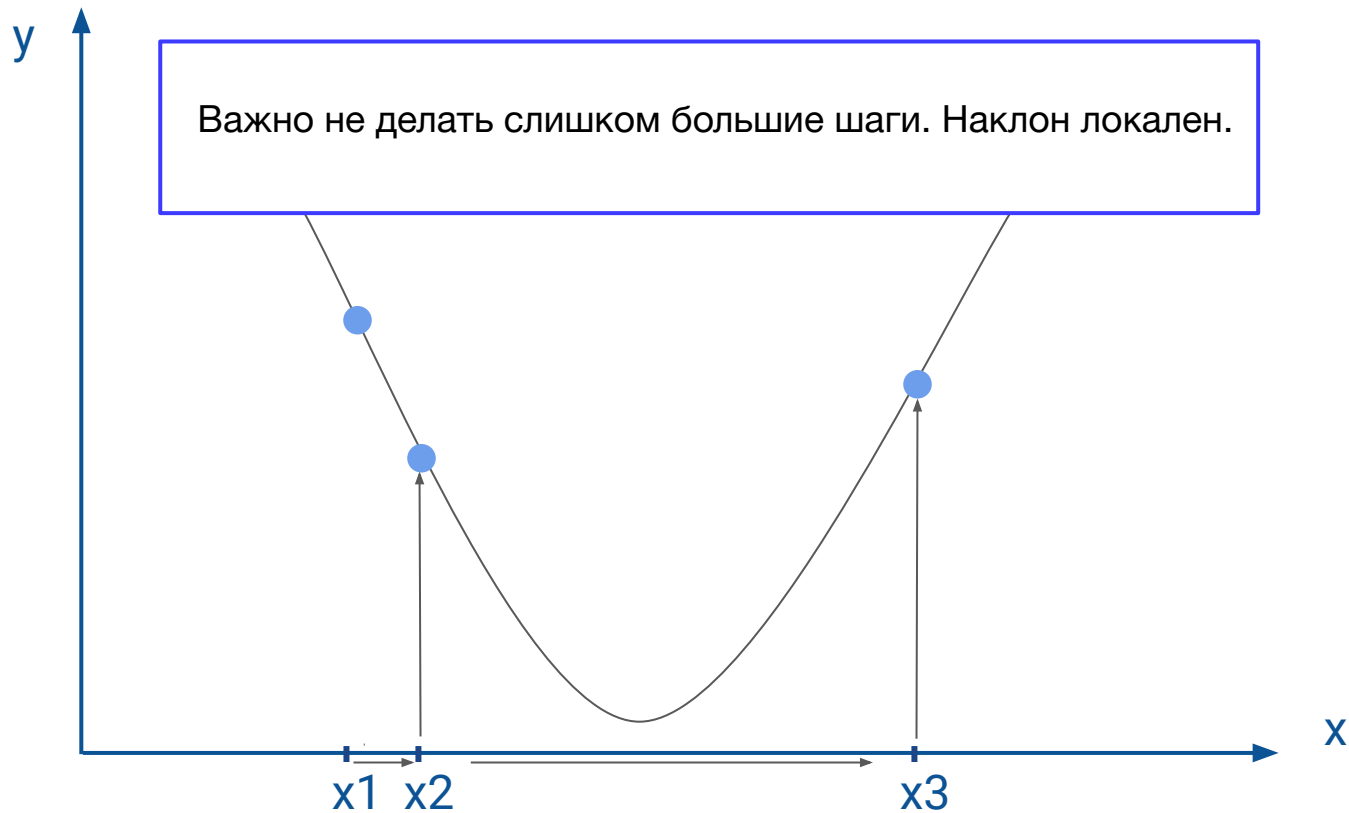
Настройка  
параметров

*\* градиент -- вектор из производных по каждому из параметров (для многомерного случая)*

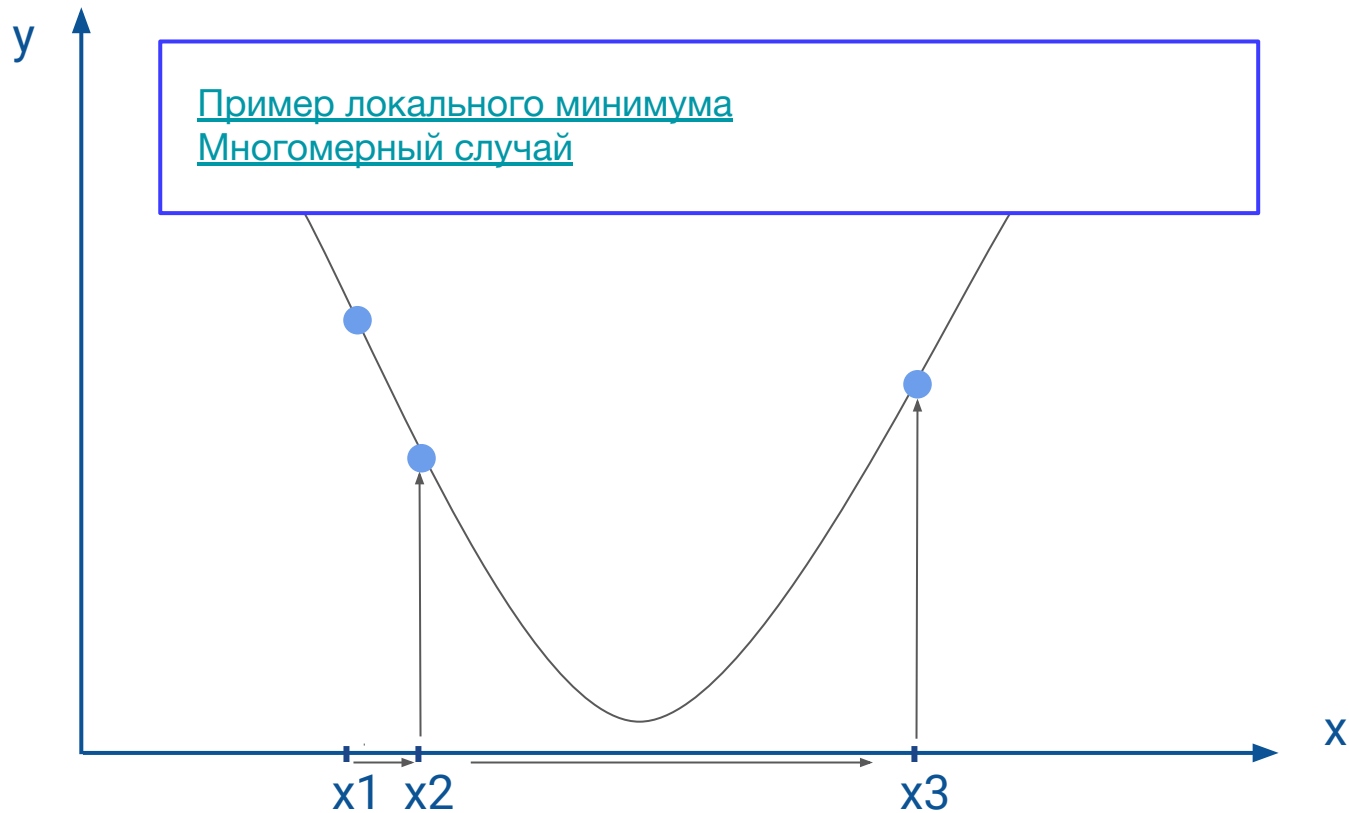
# Скорость обучения



# Скорость обучения



# Визуализация



# Градиентный спуск

- Метод прост в реализации (вы напишите его сами в практическом задании)
- Именно он (с некоторыми модификациями) используется сегодня для оптимизации нейронных сетей



# Градиентный спуск

- Метод прост в реализации (вы напишите его сами в практическом задании)
- Именно он (с некоторыми модификациями) используется сегодня для оптимизации нейронных сетей

Важные ограничения метода:

- Мы на всех этапах пользовались требованием дифференцируемости функции (без этого мы бы не могли посчитать производную)

# Градиентный спуск

- Метод прост в реализации (вы напишите его сами в практическом задании)
- Именно он (с некоторыми модификациями) используется сегодня для оптимизации нейронных сетей

Важные ограничения метода:

- Мы на всех этапах пользовались требованием дифференцируемости функции (без этого мы бы не могли посчитать производную).
- Нахождение глобального минимума не гарантировано.
- Одной из самых серьезных проблем данного подхода является застревание в локальных минимумах
- Метод очень чувствителен к инициализации

# Метод обратного распространения ошибки

# Градиентный спуск

1. Инициализировать параметры случайным образом (выбрать  $\theta_0$ )
2. Вычислить значение функционала  $\mathcal{L}$
3. Вычислить  $\frac{\partial \mathcal{L}}{\partial \theta}(\theta_{i-1})$  при текущем значении параметра

4. Сделать шаг оптимизации:

$$\theta_i = \theta_{i-1} - \lambda \frac{\partial \mathcal{L}}{\partial \theta}(\theta_{i-1})$$

**а. После этого шага параметры изменились!**

5. Повторять 2-4 пока  $\mathcal{L}$  не перестанет меняться

# Градиентный спуск

1. Инициализировать параметры случайным образом (выбрать  $\theta_0$ )
2. Вычислить значение функционала  $\mathcal{L}$

3. Вычислить  $\frac{\partial \mathcal{L}}{\partial \theta}(\theta_{i-1})$  при текущем значении параметра

4. Сделать шаг оптимизации:

$$\theta_i = \theta_{i-1} - \lambda \frac{\partial \mathcal{L}}{\partial \theta}(\theta_{i-1})$$

**а. После этого шага параметры изменились!**

5. Повторять 2-4 пока  $\mathcal{L}$  не перестанет меняться

# Вычисление производных

- Нейронная сеть: 10 -> 512 -> 256 -> 1
- Содержит около миллиона параметров
- Наивный подход:
  - Посчитать производные для каждого параметра в общем виде
  - Подставлять текущие значения параметров для вычисления производной
- Для одного шага градиентного спуска мы вычислим 1 000 000 раз значения для функций сравнимых по сложности с исходной!
- Квадратичная сложность шага

# Вычисление производных

- Нейронная сеть:  $10 \times 512 \times 256 \times 1$
- Содержит около миллиона параметров
- Наивный подход:
  - Посчитать производные каждого параметра в виде
  - Подставлять значения параметров в формулы вычисления
- Для одного шага градиента нужно вычислять миллионы раз значения для функции сравнимых по сложности с исходной!
- Квадратичная сложность шага

**Back propagation** может это делать линейно!  
С помощью Chain Rule

*\*(Back propagation -- подход к вычислению производных, обучение происходит методом градиентного спуска)*

# Chain Rule

$$f(x, y, z) = (2x + y)z$$

Найти:  $\frac{\partial f}{\partial x} \quad \frac{\partial f}{\partial y} \quad \frac{\partial f}{\partial z}$

В точке:  $x = 10, y = -2, z = -4$



# Chain Rule

$$f(x, y, z) = (2x + y)z$$

Найти:  $\frac{\partial f}{\partial x} \quad \frac{\partial f}{\partial y} \quad \frac{\partial f}{\partial z}$

В точке:  $x = 10, y = -2, z = -4$

$$\frac{\partial f}{\partial x} = 2z$$

$$\frac{\partial f}{\partial y} = z$$

$$\frac{\partial f}{\partial z} = 2x + y$$

# Chain Rule

$$f(x, y, z) = (2x + y)z$$

Найти:  $\frac{\partial f}{\partial x} \quad \frac{\partial f}{\partial y} \quad \frac{\partial f}{\partial z}$

В точке:  $x = 10, y = -2, z = -4$

$$\frac{\partial f}{\partial x} = 2z = -8$$

$$\frac{\partial f}{\partial y} = z = -4$$

$$\frac{\partial f}{\partial z} = 2x + y = 18$$

# Chain Rule

$$f(x, y, z) = (2x + y)z$$

Найти:  $\frac{\partial f}{\partial x} \quad \frac{\partial f}{\partial y} \quad \frac{\partial f}{\partial z}$

В точке:  $x = 10, y = -2, z = -4$

$$\frac{\partial f}{\partial x} = \boxed{2}z = -8$$

$$\frac{\partial f}{\partial y} = z = -4$$

$$\frac{\partial f}{\partial z} = 2x + y = 18$$

# Chain Rule

$$f(x, y, z) = (2x + y)z$$

Найти:  $\frac{\partial f}{\partial x} \quad \frac{\partial f}{\partial y} \quad \frac{\partial f}{\partial z}$

В точке:  $x = 10, y = -2, z = -4$

$$\frac{\partial f}{\partial x} = 2z = -8$$

$$\frac{\partial f}{\partial y} = z = -4$$

$$\frac{\partial f}{\partial z} = 2x + y = 18$$

$$\frac{\partial f(g(x))}{\partial x} = \frac{\partial f}{\partial g} \frac{\partial g}{\partial x}$$

# Chain Rule

$$f(x, y, z) = (2x + y)z$$

Найти:  $\frac{\partial f}{\partial x} \quad \frac{\partial f}{\partial y} \quad \frac{\partial f}{\partial z}$

В точке:  $x = 10, y = -2, z = -4$

$$\frac{\partial f}{\partial x} = 2z = -8$$

$$\frac{\partial f}{\partial y} = z = -4$$

$$\frac{\partial f}{\partial z} = 2x + y = 18$$

$$\frac{\partial f(g(x))}{\partial x} = \frac{\partial f}{\partial g} \frac{\partial g}{\partial x}$$

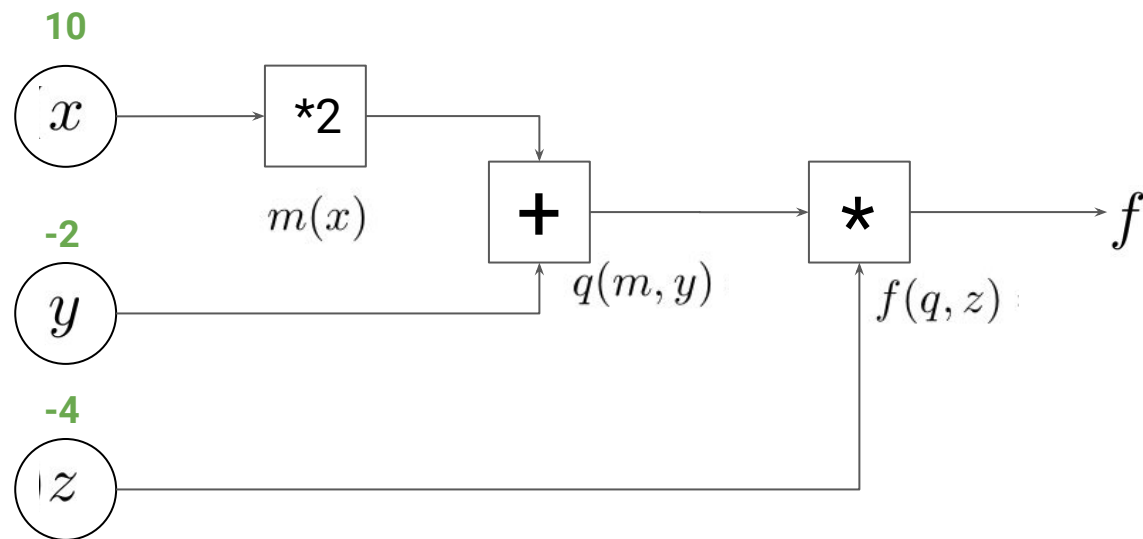
$$f(g, z) = gz$$

$$g(x, y) = 2x + y$$

# Вычислительный граф

$$f(x, y, z) = (2x + y)z$$

$$x = 10, y = -2, z = -4$$



$$m(x) = 2x$$

$$q(m, y) = m + y$$

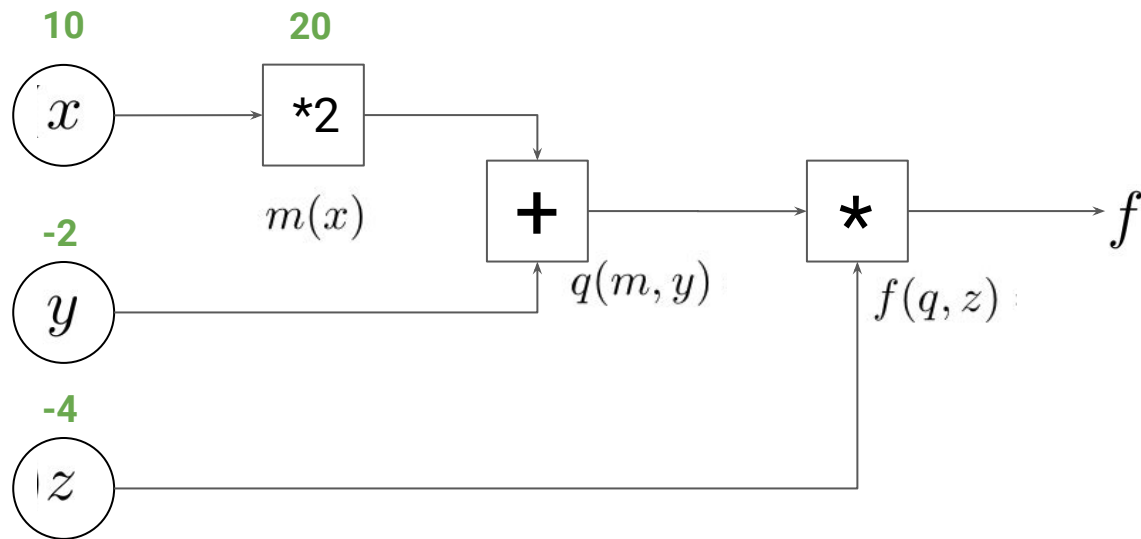
$$f(q, z) = qz$$

# Вычислительный граф

Forward pass

$$f(x, y, z) = (2x + y)z$$

$$x = 10, y = -2, z = -4$$



$$m(x) = 2x$$

$$q(m, y) = m + y$$

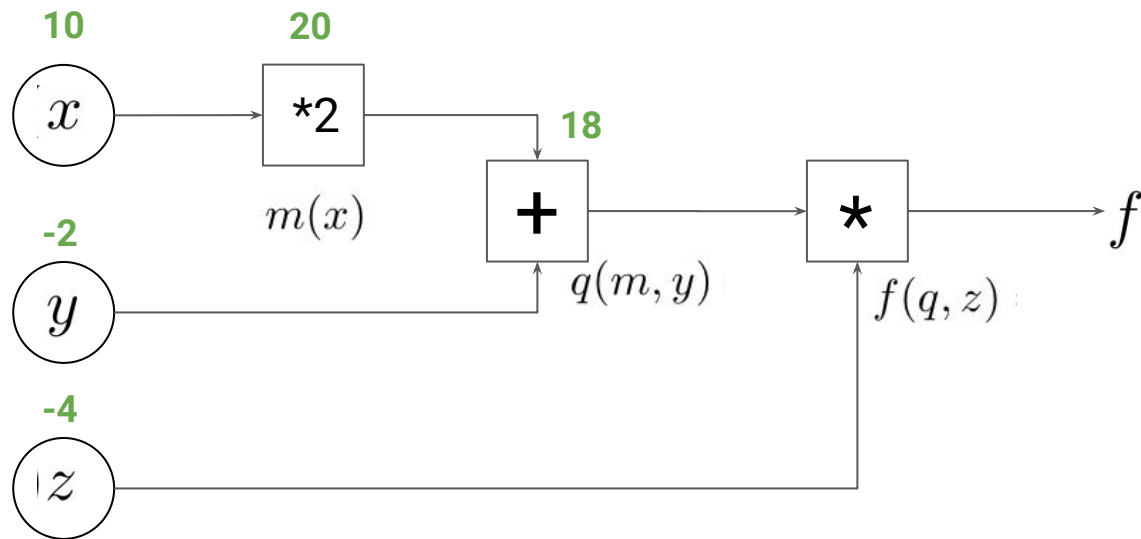
$$f(q, z) = qz$$

# Вычислительный граф

Forward pass

$$f(x, y, z) = (2x + y)z$$

$$x = 10, y = -2, z = -4$$



$$m(x) = 2x$$

$$q(m, y) = m + y$$

$$f(q, z) = qz$$

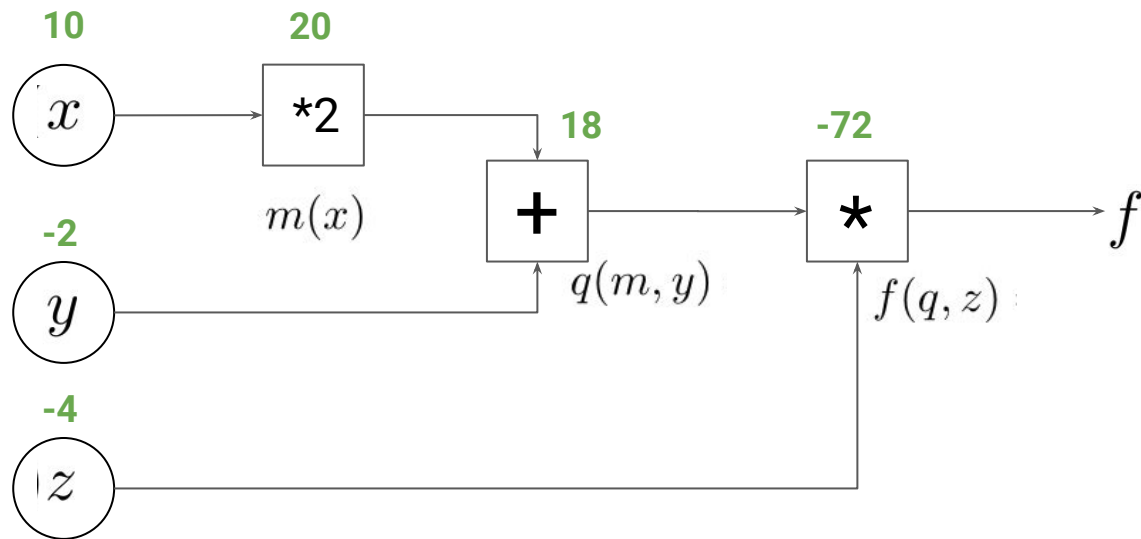


# Вычислительный граф

Forward pass

$$f(x, y, z) = (2x + y)z$$

$$x = 10, y = -2, z = -4$$



$$m(x) = 2x$$

$$q(m, y) = m + y$$

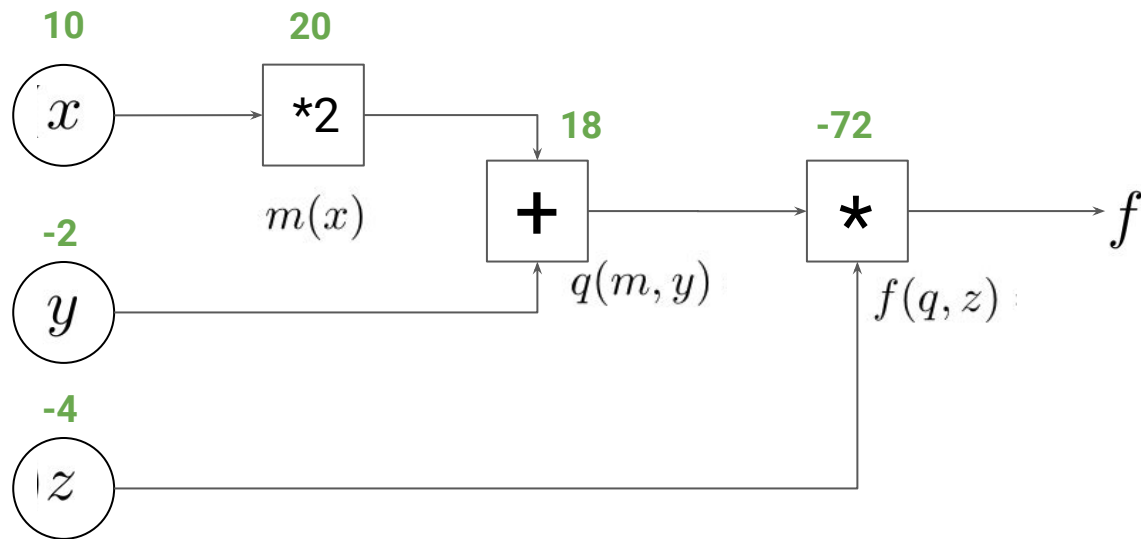
$$f(q, z) = qz$$

# Вычислительный граф

Backward pass

$$f(x, y, z) = (2x + y)z$$

$$x = 10, y = -2, z = -4$$



$$m(x) = 2x$$

$$q(m, y) = m + y$$

$$f(q, z) = qz$$

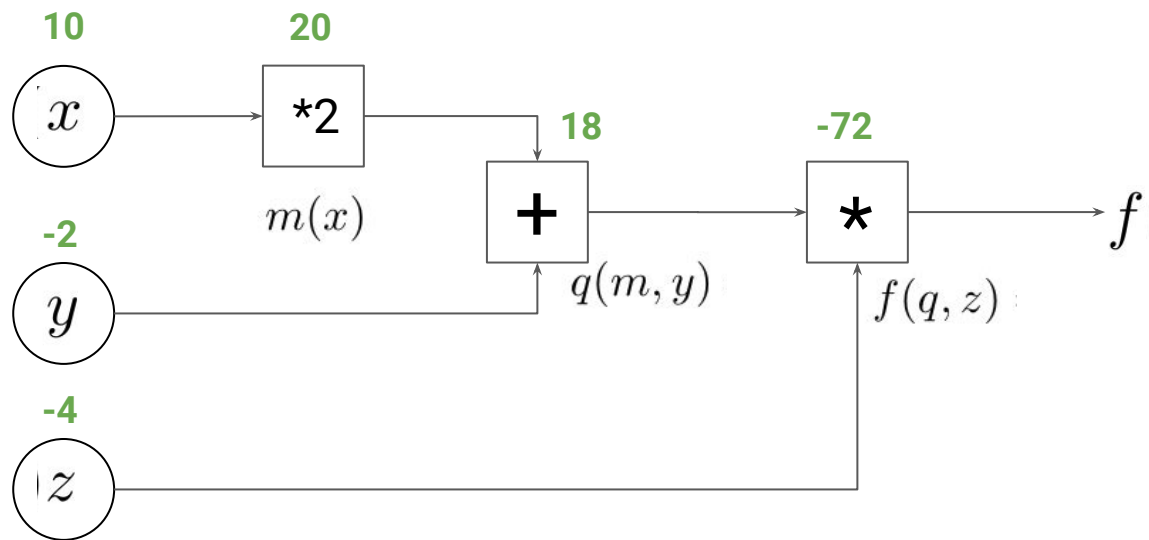
$$f(x, y, z) = (2x + y)z$$

$$x = 10, y = -2, z = -4$$

$$m(x) = 2x$$

$$q(m, y) = m + y$$

$$f(q, z) = qz$$



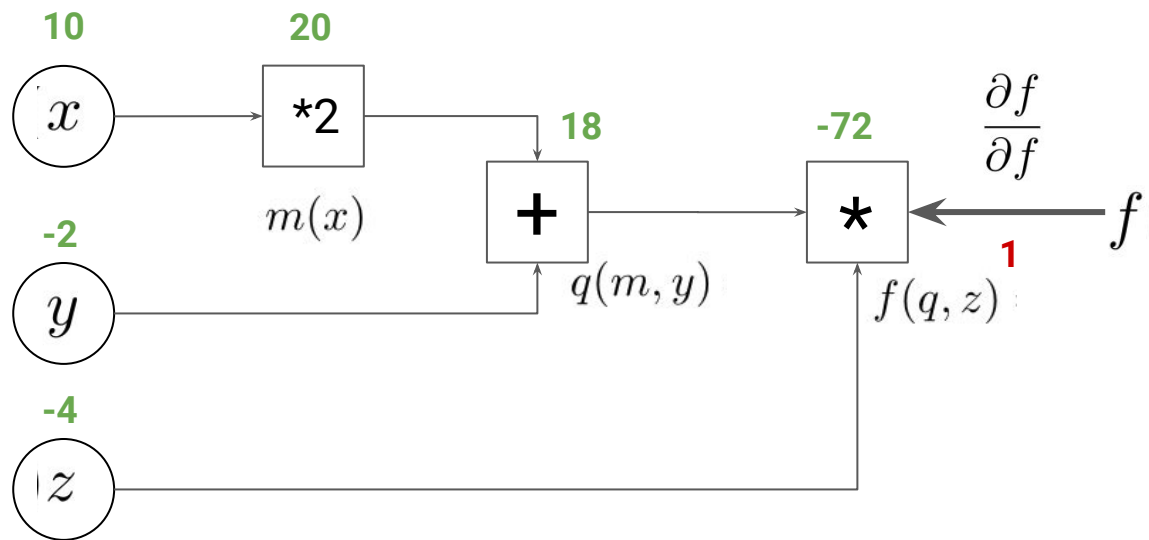
$$f(x, y, z) = (2x + y)z$$

$$x = 10, y = -2, z = -4$$

$$m(x) = 2x$$

$$q(m, y) = m + y$$

$$f(q, z) = qz$$



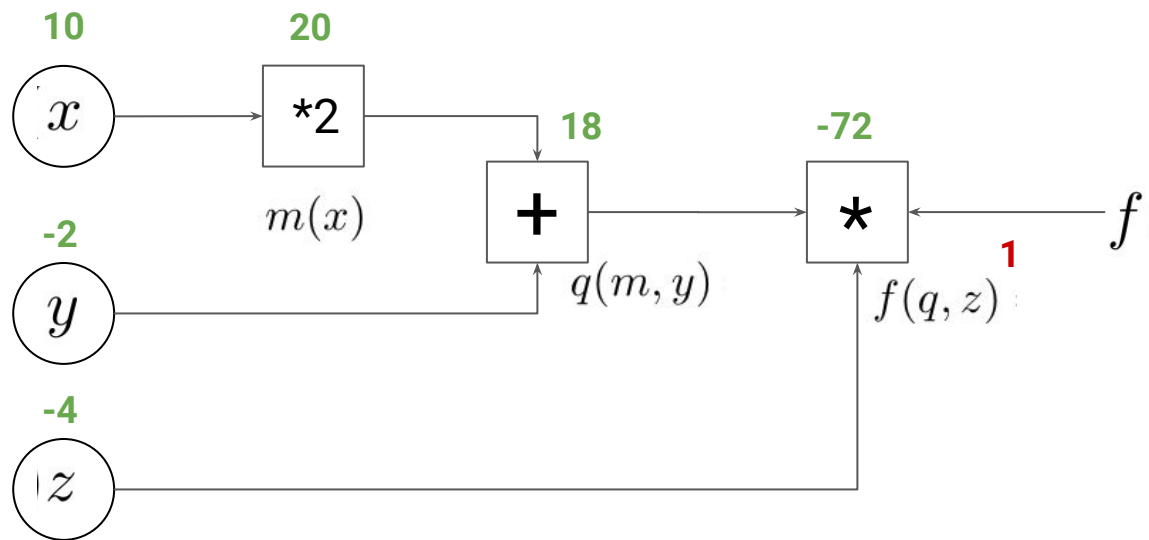
$$f(x, y, z) = (2x + y)z$$

$$x = 10, y = -2, z = -4$$

$$m(x) = 2x$$

$$q(m, y) = m + y$$

$$f(q, z) = qz$$



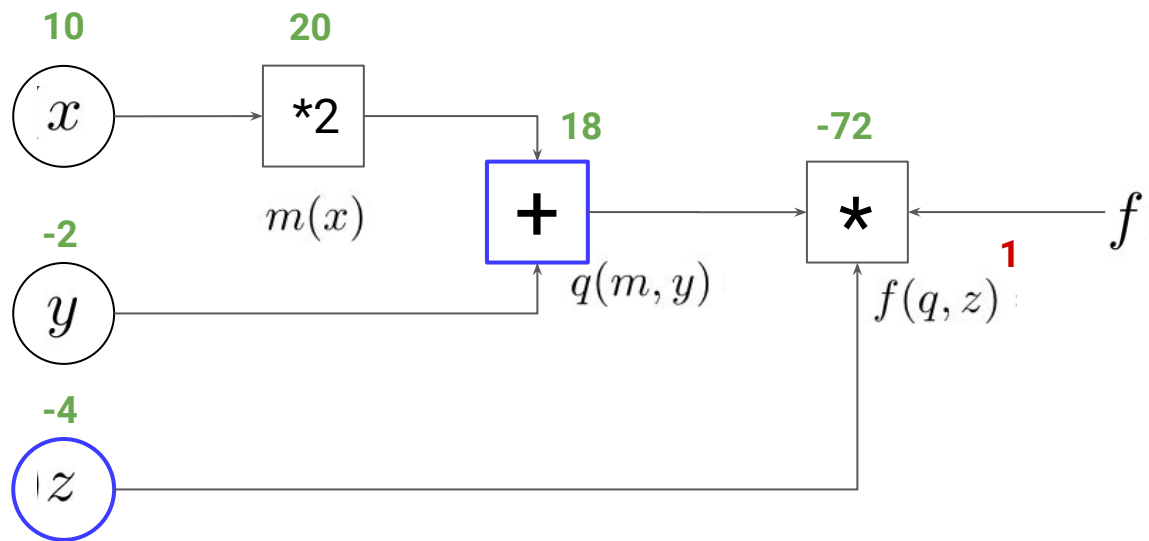
$$f(x, y, z) = (2x + y)z$$

$$x = 10, y = -2, z = -4$$

$$m(x) = 2x$$

$$q(m, y) = m + y$$

$$f(q, z) = qz$$



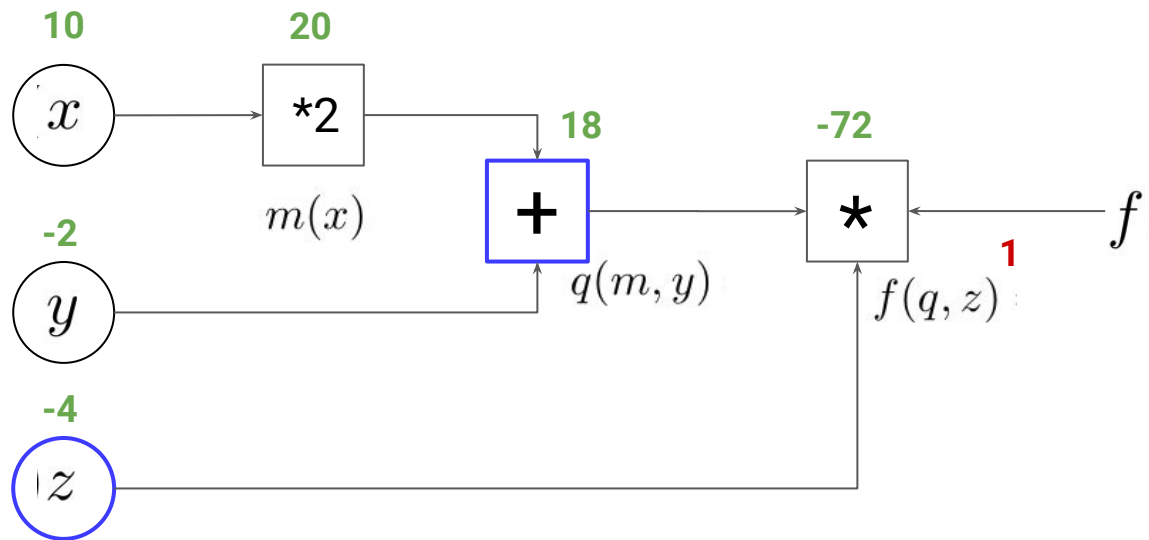
$$f(x, y, z) = (2x + y)z$$

$$x = 10, y = -2, z = -4$$

$$m(x) = 2x$$

$$q(m, y) = m + y$$

$$f(q, z) = qz$$



$$\frac{\partial f}{\partial z} = q$$

$$\frac{\partial f}{\partial q} = z$$

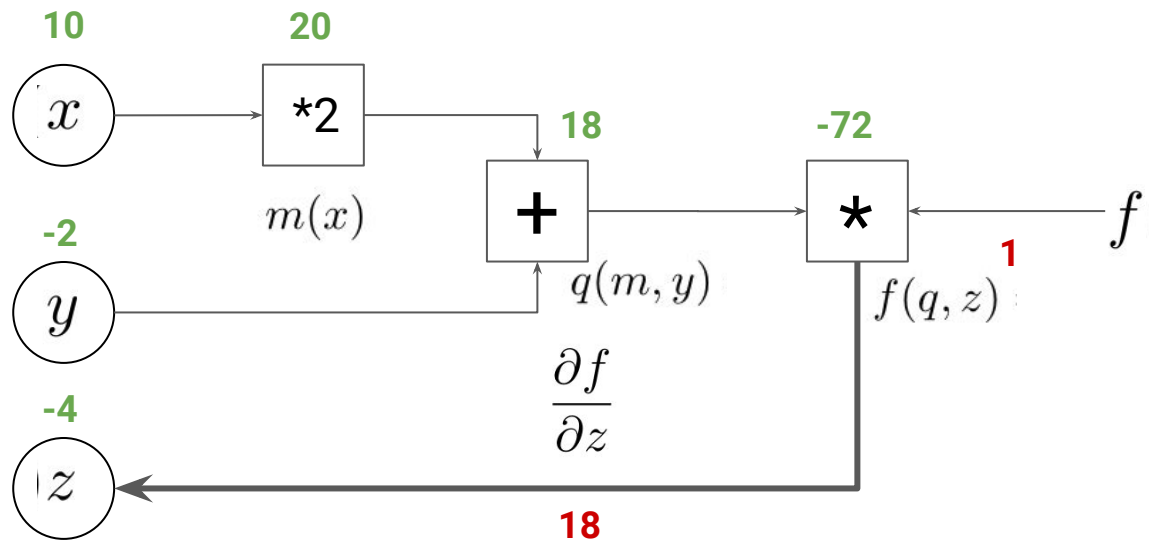
$$f(x, y, z) = (2x + y)z$$

$$x = 10, y = -2, z = -4$$

$$m(x) = 2x$$

$$q(m, y) = m + y$$

$$f(q, z) = qz$$



$$\frac{\partial f}{\partial z} = q$$

$$\frac{\partial f}{\partial q} = z$$



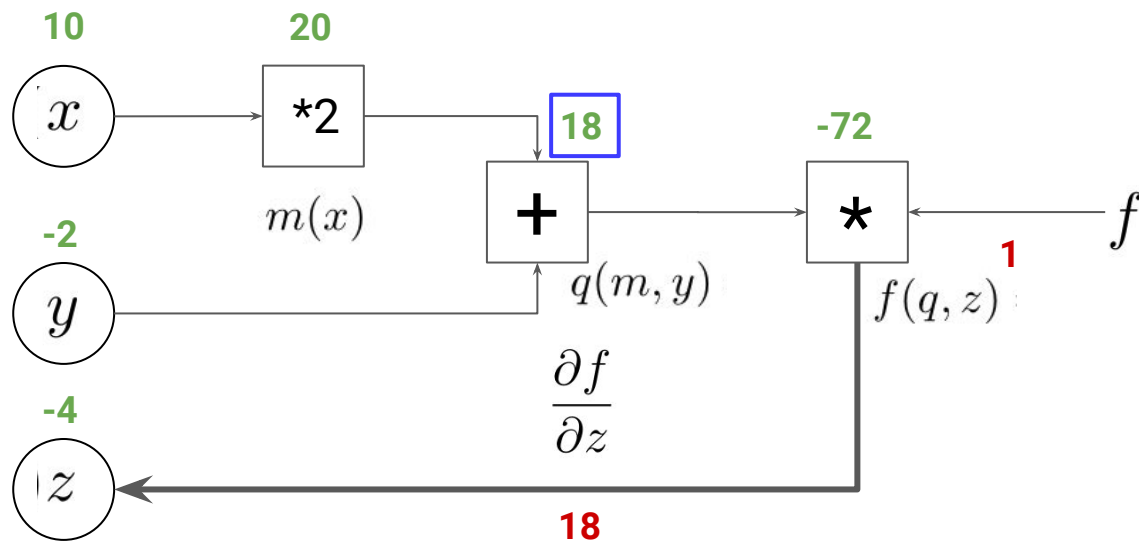
$$f(x, y, z) = (2x + y)z$$

$$x = 10, y = -2, z = -4$$

$$m(x) = 2x$$

$$q(m, y) = m + y$$

$$f(q, z) = qz$$



$$\frac{\partial f}{\partial z} = q$$

$$\frac{\partial f}{\partial q} = z$$

Мы не считали это значение!  
Оно нам известно!

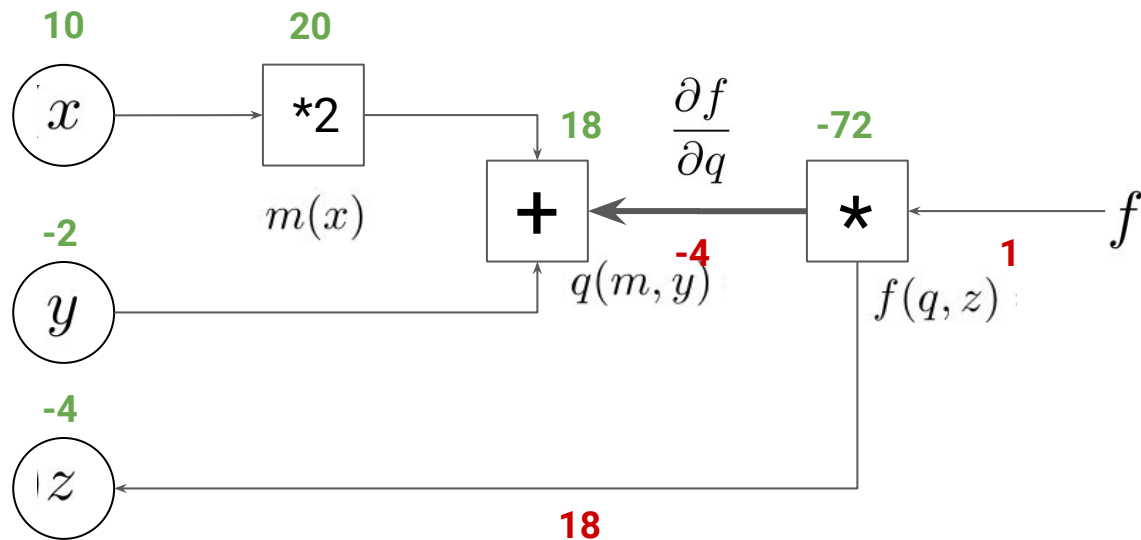
$$f(x, y, z) = (2x + y)z$$

$$x = 10, y = -2, z = -4$$

$$m(x) = 2x$$

$$q(m, y) = m + y$$

$$f(q, z) = qz$$



$$\frac{\partial f}{\partial z} = q$$

$$\frac{\partial f}{\partial q} = z$$

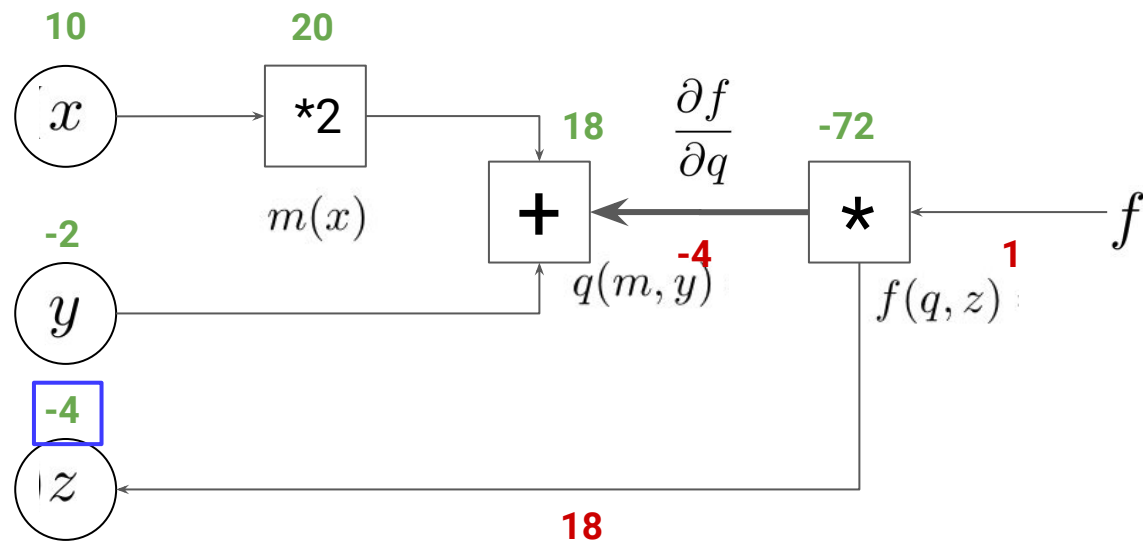
$$f(x, y, z) = (2x + y)z$$

$$x = 10, y = -2, z = -4$$

$$m(x) = 2x$$

$$q(m, y) = m + y$$

$$f(q, z) = qz$$



$$\frac{\partial f}{\partial z} = q$$

$$\frac{\partial f}{\partial q} = z$$

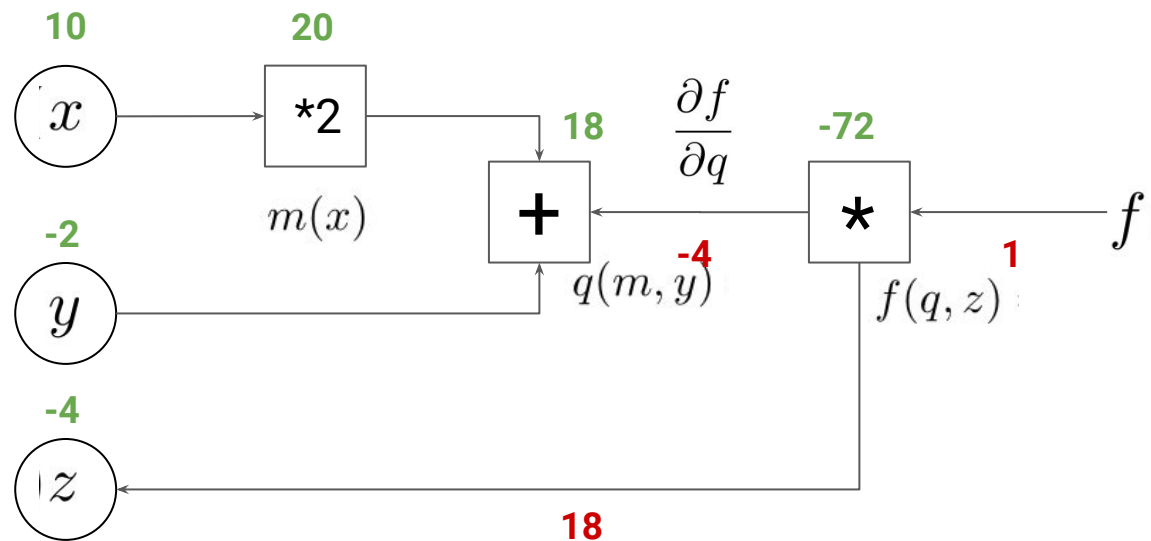
$$f(x, y, z) = (2x + y)z$$

$$x = 10, y = -2, z = -4$$

$$m(x) = 2x$$

$$q(m, y) = m + y$$

$$f(q, z) = qz$$



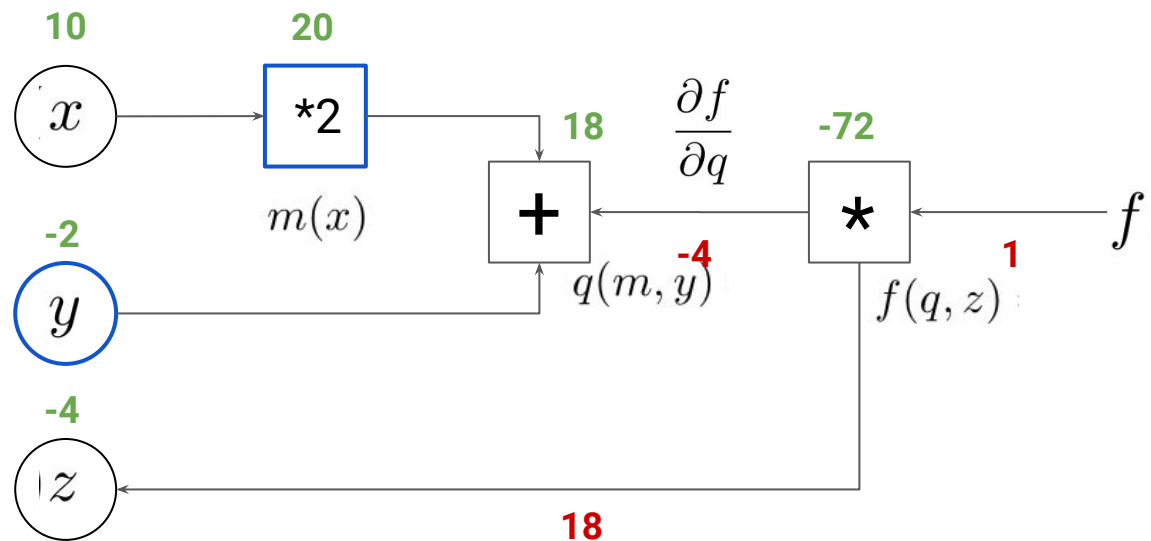
$$f(x, y, z) = (2x + y)z$$

$$x = 10, y = -2, z = -4$$

$$m(x) = 2x$$

$$q(m, y) = m + y$$

$$f(q, z) = qz$$



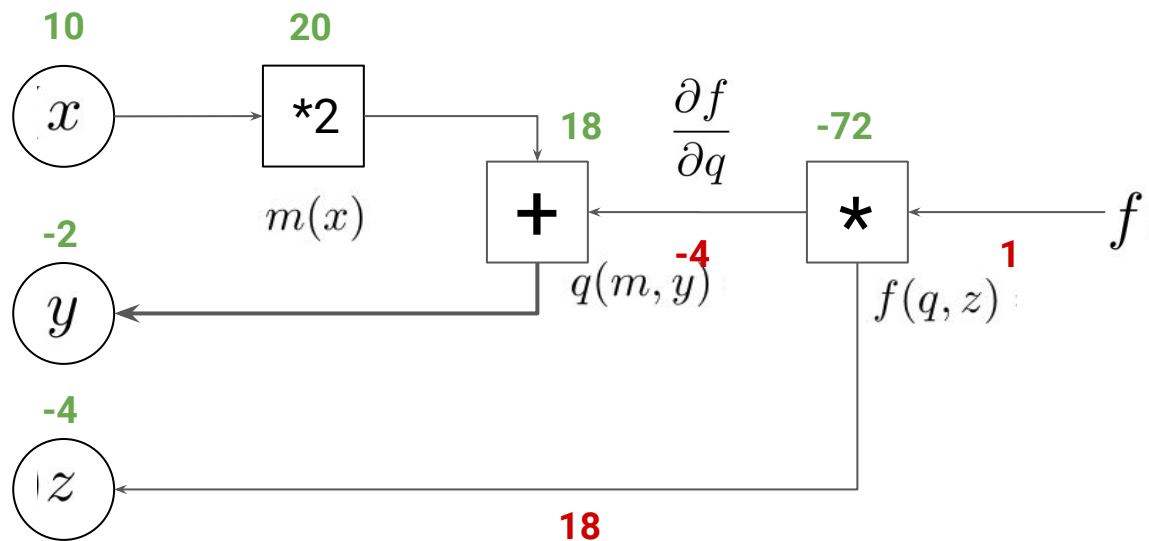
$$f(x, y, z) = (2x + y)z$$

$$x = 10, y = -2, z = -4$$

$$m(x) = 2x$$

$$q(m, y) = m + y$$

$$f(q, z) = qz$$



$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial y}$$

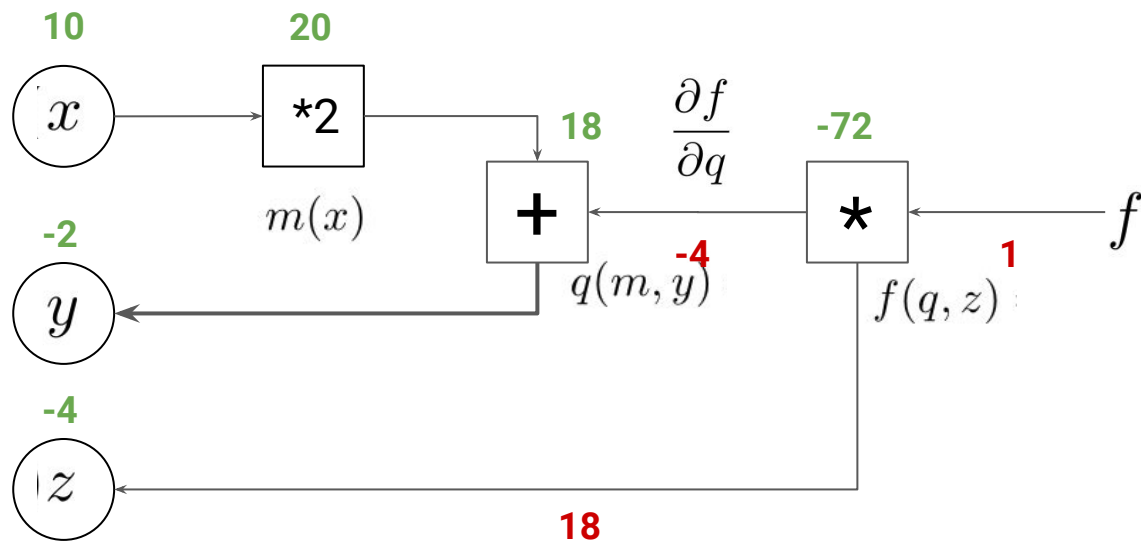
$$f(x, y, z) = (2x + y)z$$

$$x = 10, y = -2, z = -4$$

$$m(x) = 2x$$

$$q(m, y) = m + y$$

$$f(q, z) = qz$$



$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial y}$$

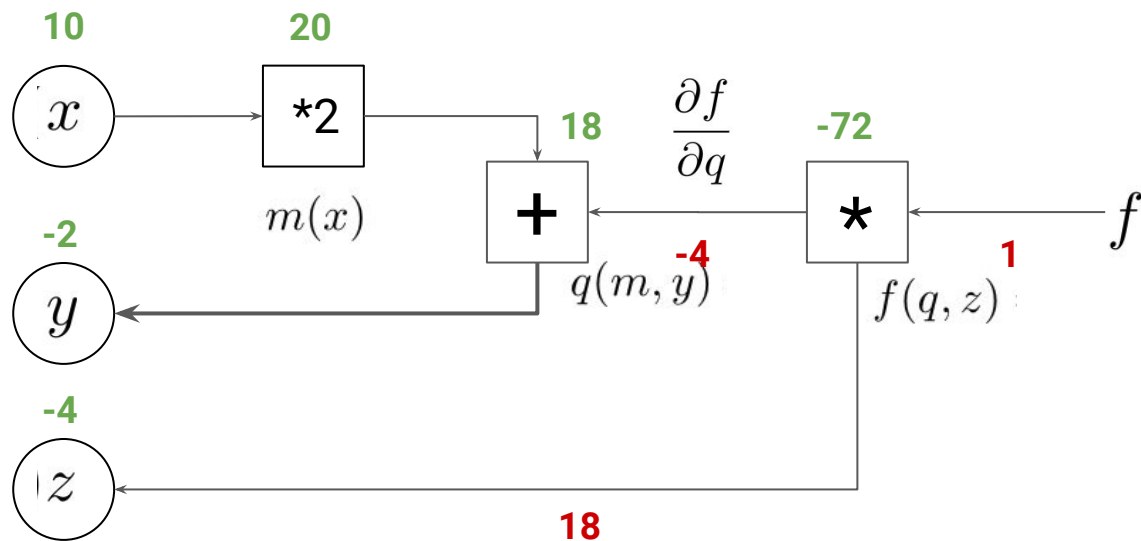
$$f(x, y, z) = (2x + y)z$$

$$x = 10, y = -2, z = -4$$

$$m(x) = 2x$$

$$q(m, y) = m + y$$

$$f(q, z) = qz$$



$$\frac{\partial f}{\partial y} = \boxed{\frac{\partial f}{\partial q}} \frac{\partial q}{\partial y}$$



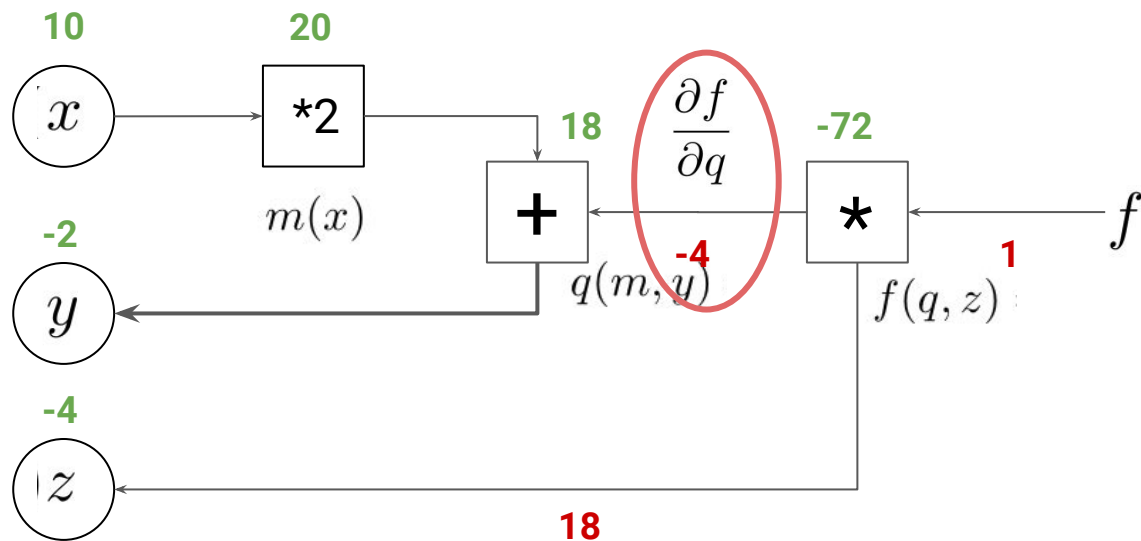
$$f(x, y, z) = (2x + y)z$$

$$x = 10, y = -2, z = -4$$

$$m(x) = 2x$$

$$q(m, y) = m + y$$

$$f(q, z) = qz$$



$$\frac{\partial f}{\partial y} = \boxed{\frac{\partial f}{\partial q}} \frac{\partial q}{\partial y}$$

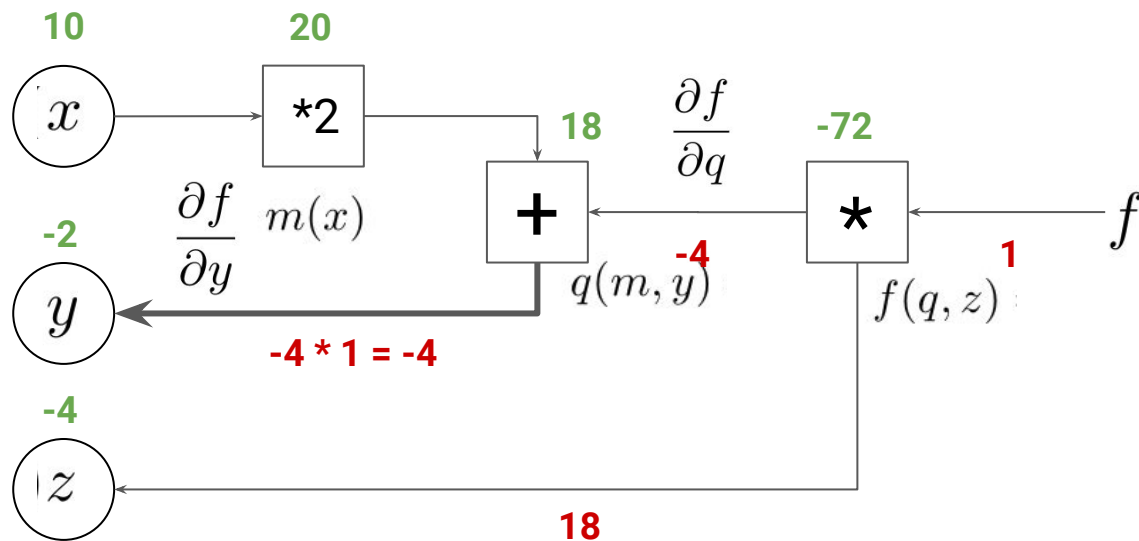
$$f(x, y, z) = (2x + y)z$$

$$x = 10, y = -2, z = -4$$

$$m(x) = 2x$$

$$q(m, y) = m + y$$

$$f(q, z) = qz$$



$$\frac{\partial f}{\partial y} = \boxed{\frac{\partial f}{\partial q}} \frac{\partial q}{\partial y}$$

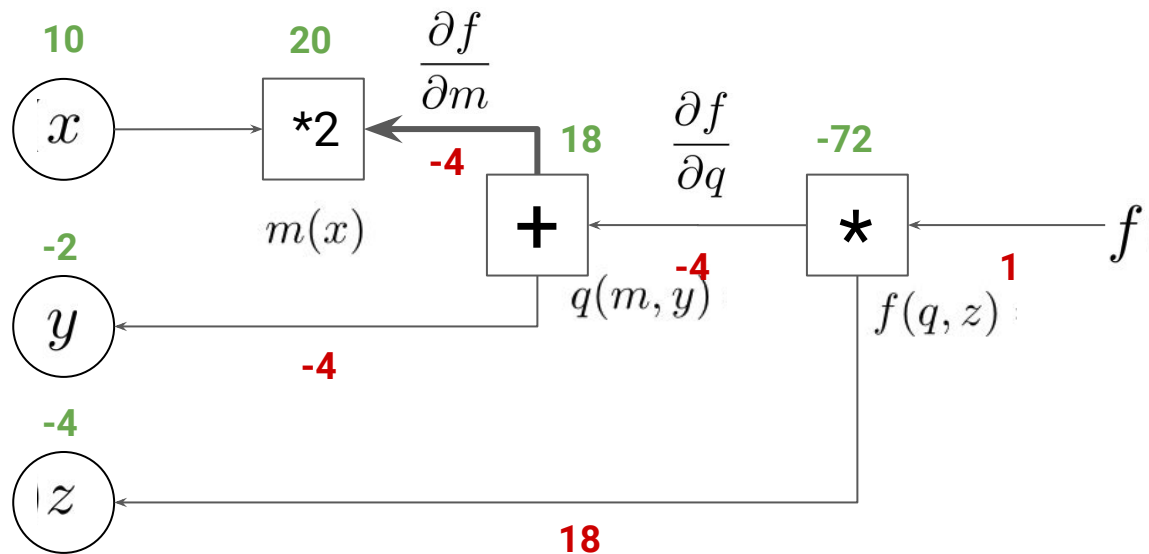
$$f(x, y, z) = (2x + y)z$$

$$x = 10, y = -2, z = -4$$

$$m(x) = 2x$$

$$q(m, y) = m + y$$

$$f(q, z) = qz$$



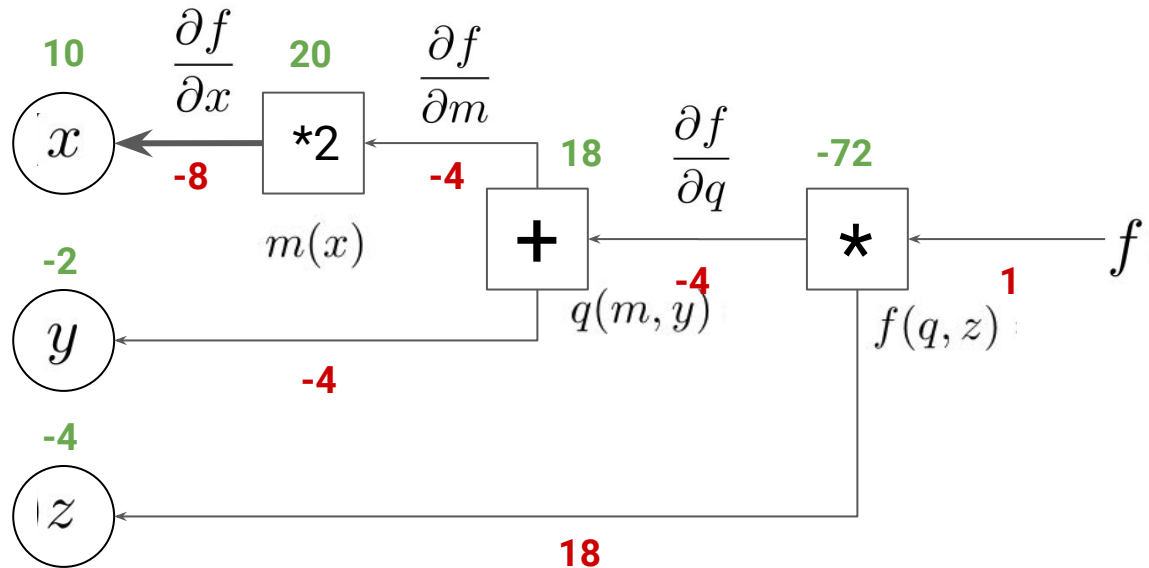
$$f(x, y, z) = (2x + y)z$$

$$x = 10, y = -2, z = -4$$

$$m(x) = 2x$$

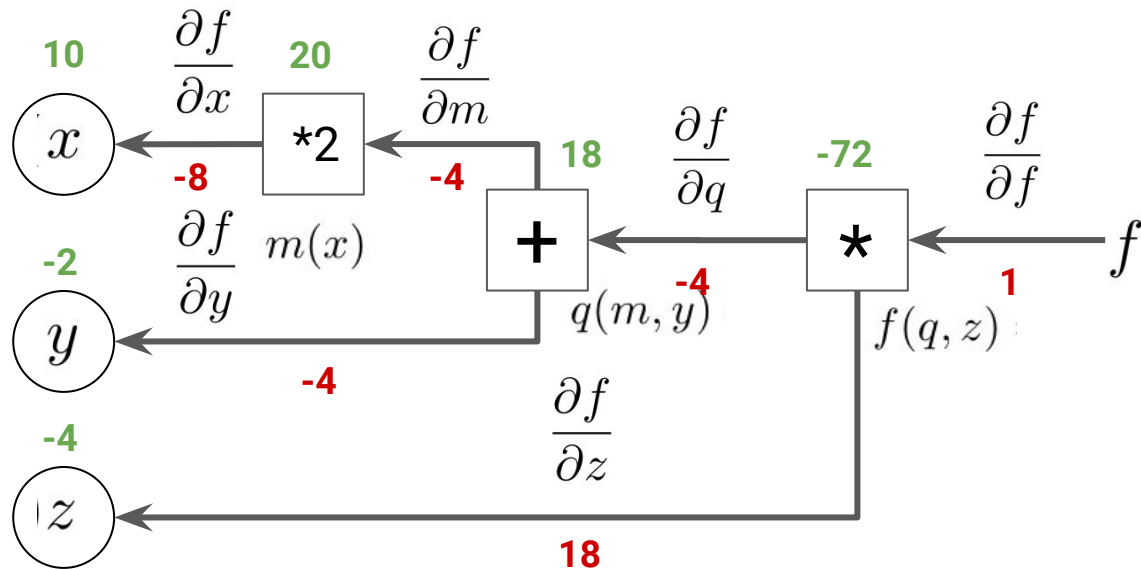
$$q(m, y) = m + y$$

$$f(q, z) = qz$$



$$f(x, y, z) = (2x + y)z$$

$$x = 10, y = -2, z = -4$$

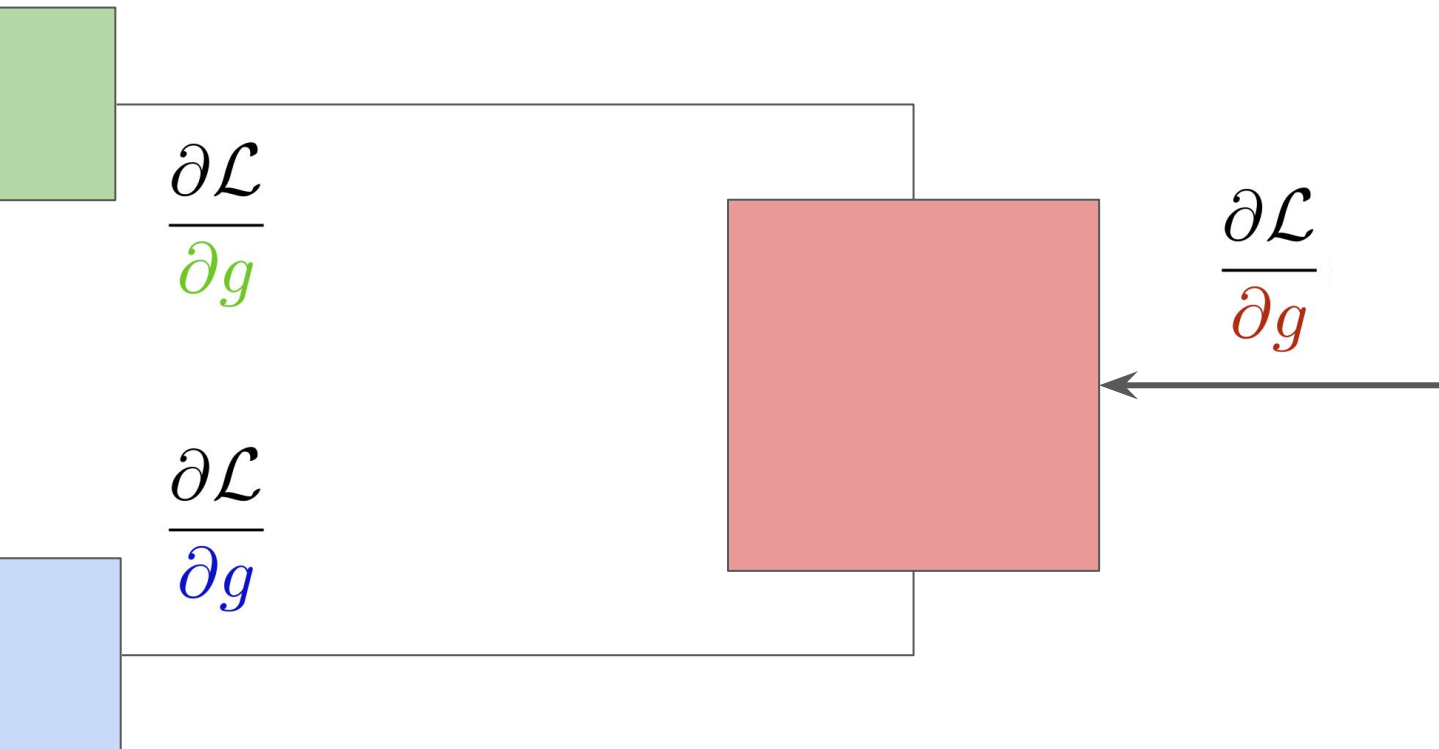


$$\frac{\partial f}{\partial x} = 2z = -8$$

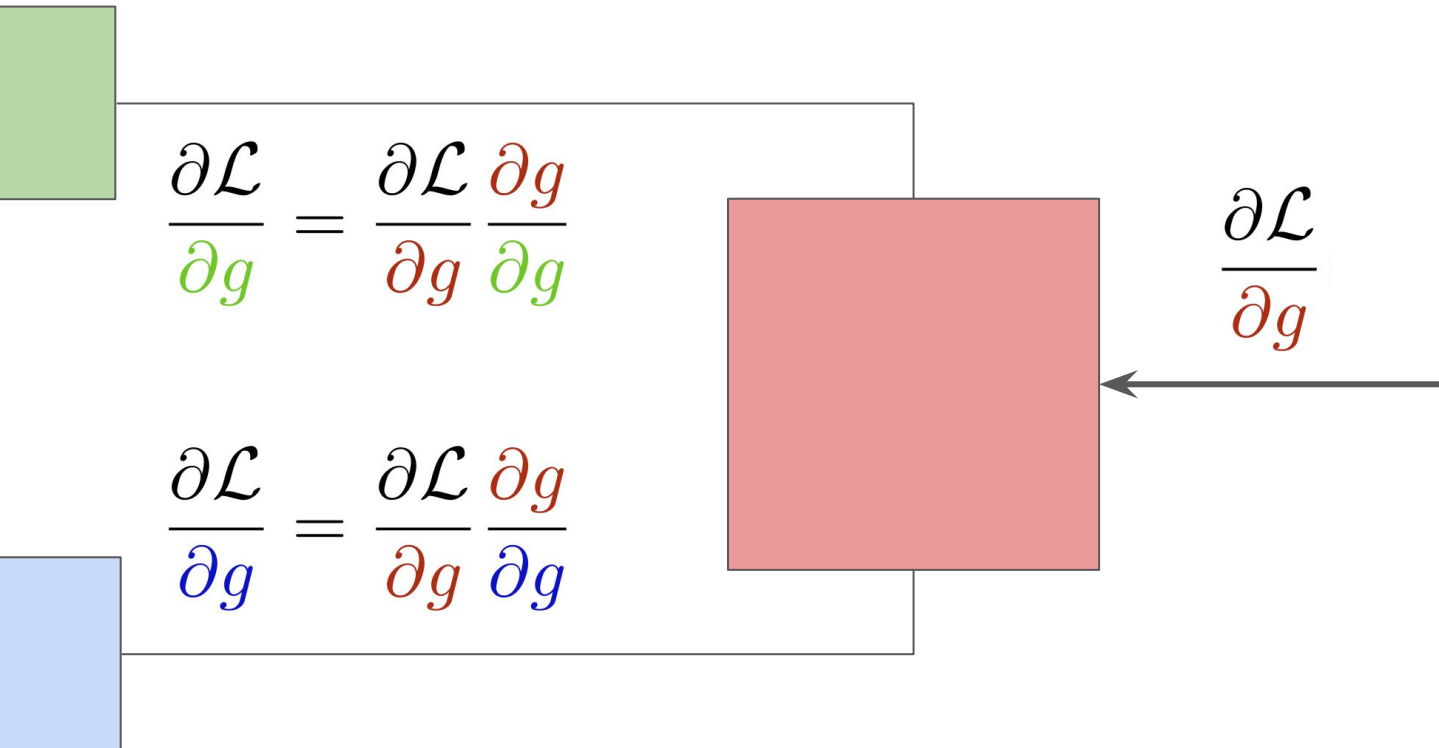
$$\frac{\partial f}{\partial y} = z = -4$$

$$\frac{\partial f}{\partial z} = 2x + y = 18$$

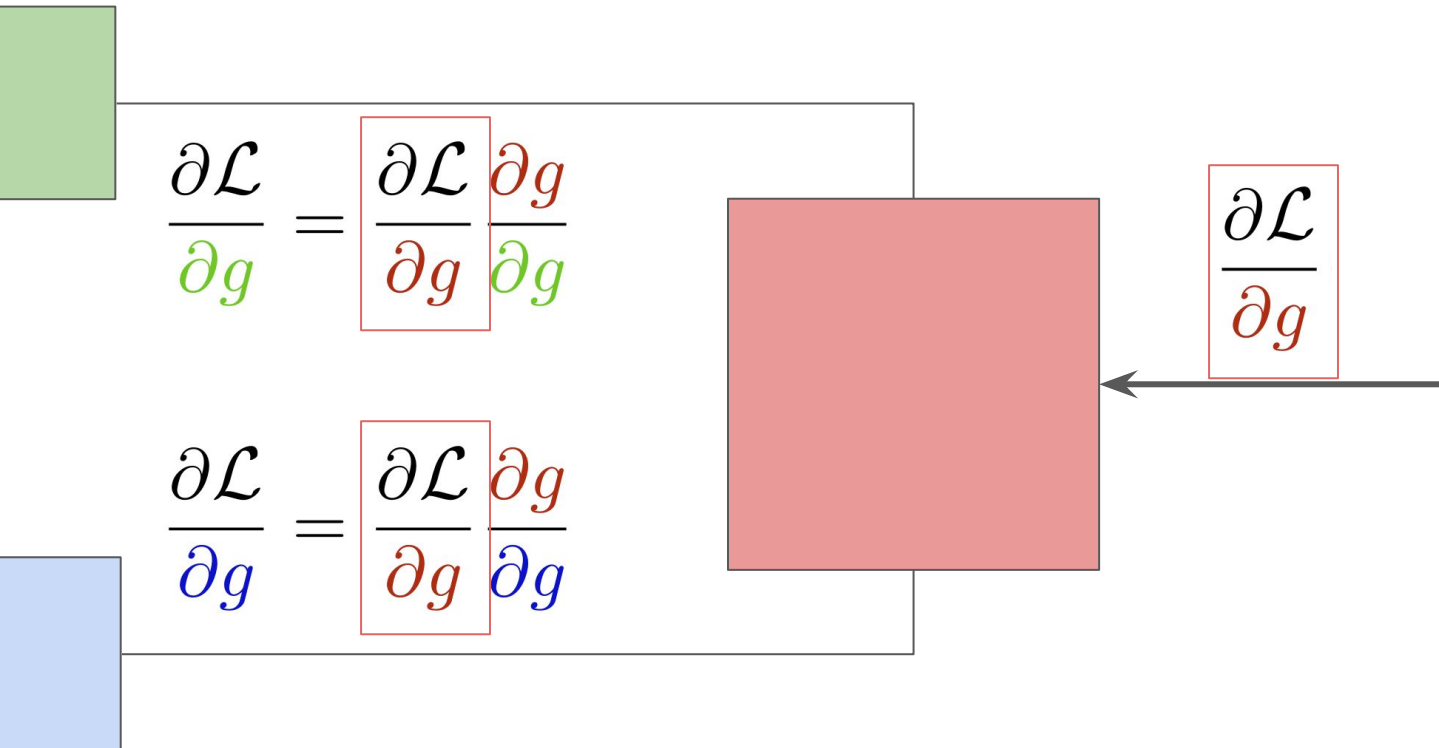
# Что должен уметь гейт



# Что должен уметь гейт

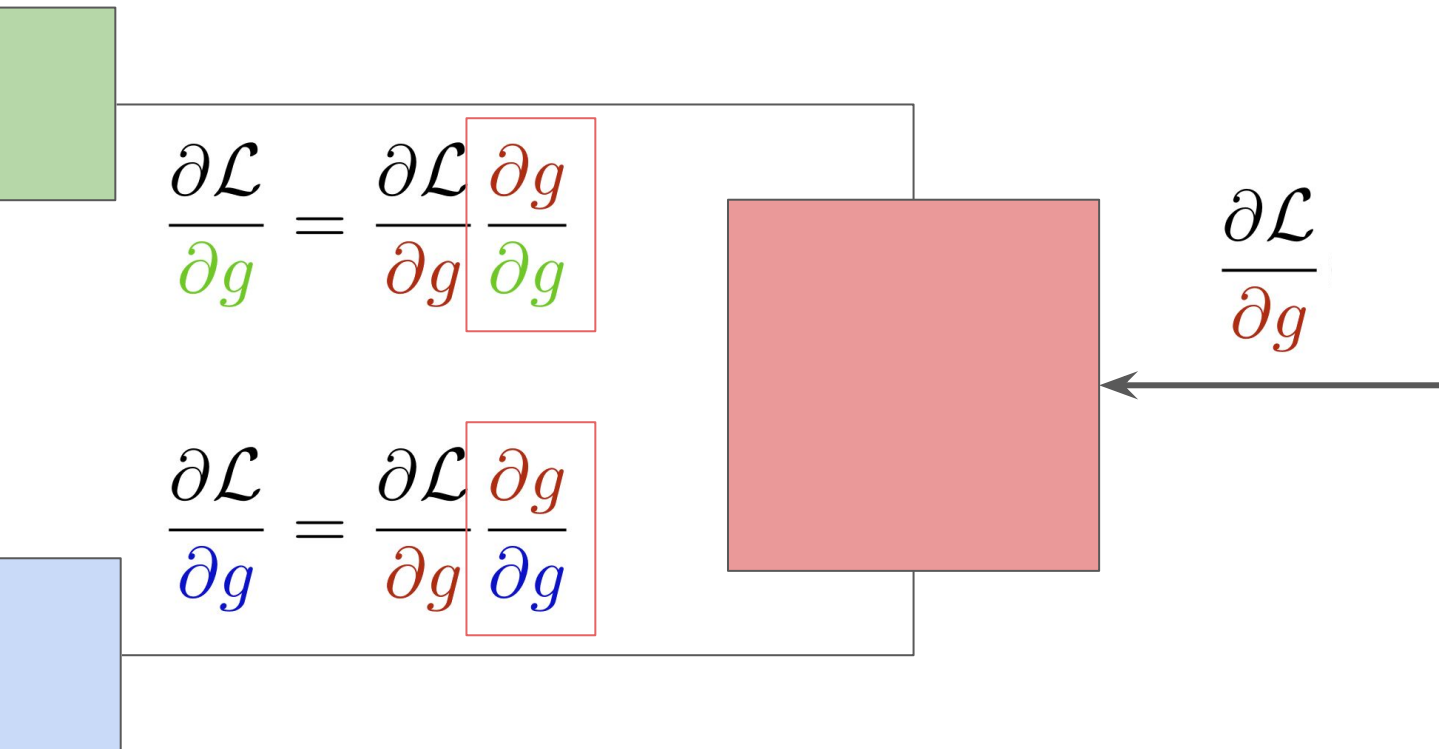


# Что должен уметь гейт





# Что должен уметь гейт

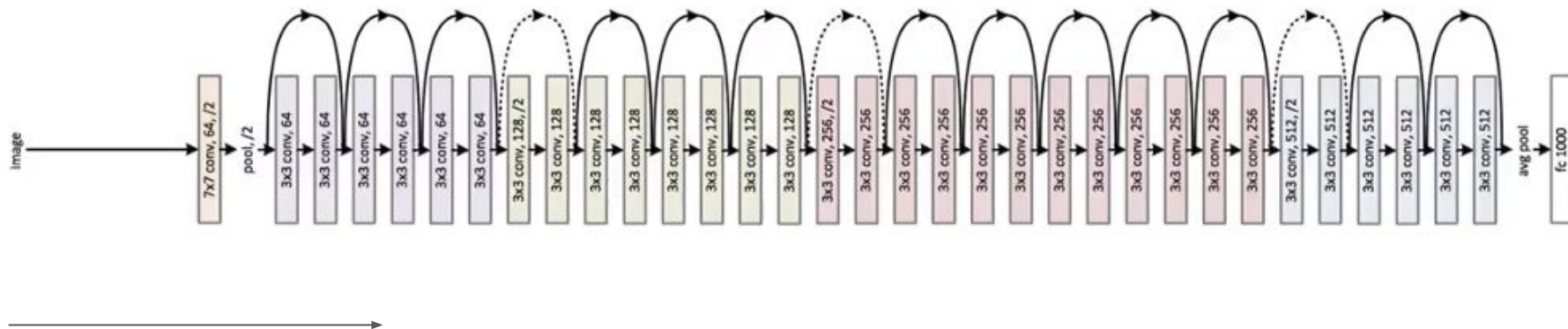


```
class MultGate(Layer):  
    def __init__(self):  
        self.x = None  
        self.y = None  
  
    def forward(self, inputs):  
        self.x, self.y = inputs  
        z = self.x * self.y  
        return z  
  
    def backward(self, dL_dz):  
        dz_dx = self.y  
        dz_dy = self.x  
  
        dL_dx = dL_dz * dz_dx  
        dL_dy = dL_dz * dz_dy  
  
        return dL_dx, dL_dy
```

```
class MultGate(Layer):  
    def __init__(self):  
        self.x = None  
        self.y = None  
  
    def forward(self, inputs):  
        self.x, self.y = inputs  
        z = self.x * self.y  
        return z  
  
    def backward(self, dL_dz):  
        dz_dx = self.y  
        dz_dy = self.x  
  
        dL_dx = dL_dz * dz_dx  
        dL_dy = dL_dz * dz_dy  
  
        return dL_dx, dL_dy
```

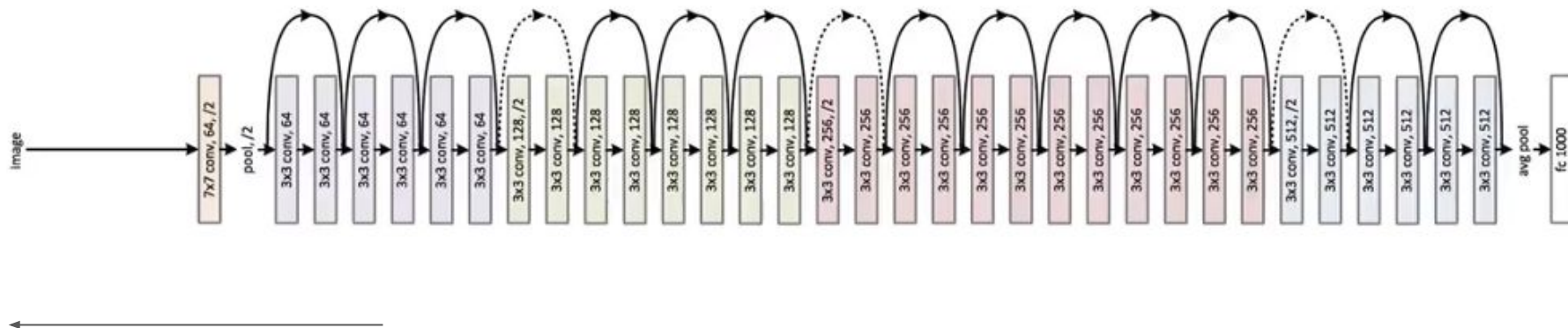
```
class MaxGate(Layer):  
    def __init__(self):  
        self.x = None  
        self.y = None  
  
    def forward(self, inputs):  
        self.x, self.y = inputs  
  
        return max(self.x, self.y)  
  
    def backward(self, dL_dz):  
  
        if self.x > self.y:  
            dz_dx = 1  
            dz_dy = 0  
        else:  
            dz_dx = 0  
            dz_dy = 1  
  
        dL_dx = dL_dz * dz_dx  
        dL_dy = dL_dz * dz_dy  
  
        return dL_dx, dL_dy
```

# Нейронная сеть как вычислительный граф



Forward каждого слоя по-очередности

# Нейронная сеть как вычислительный граф



Backward каждого слоя в обратном порядке

# Итог

- Функции можно представлять в виде вычислительного графа
- Forward pass -- считает значение функции
- Backward pass -- производные
- Back propagation позволяет быстро и эффективно считать производные сложных функций
- Эффективность достигается за счет того, что мы используем часть информации полученной в forward pass
- Для реализации гейта (слоя) нужно определить два метода: forward и backward



*On backpropagation:*

“My view is throw it all away and start again.”

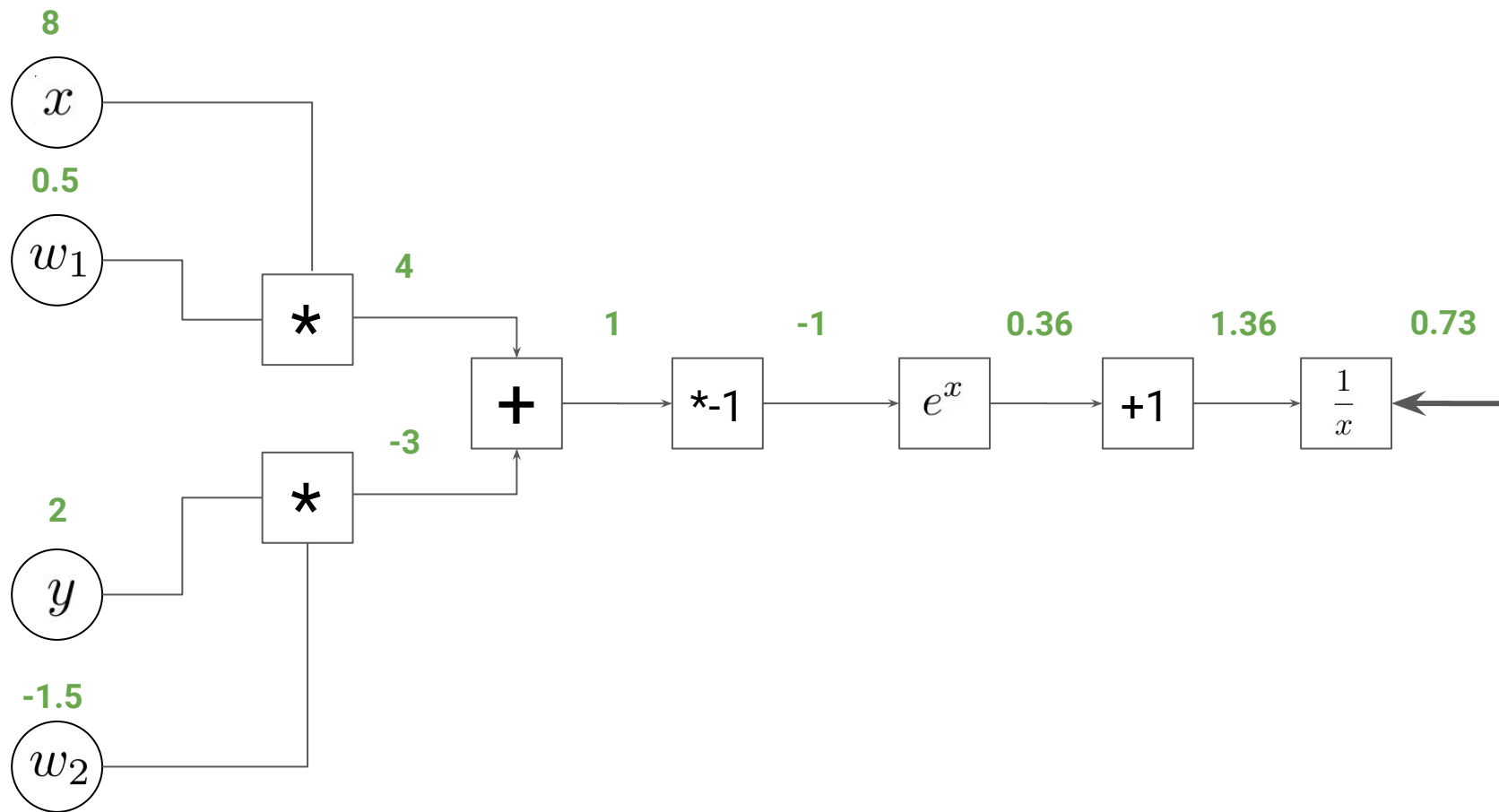
“The future depends on some graduate student who is deeply suspicious of everything I have said.”

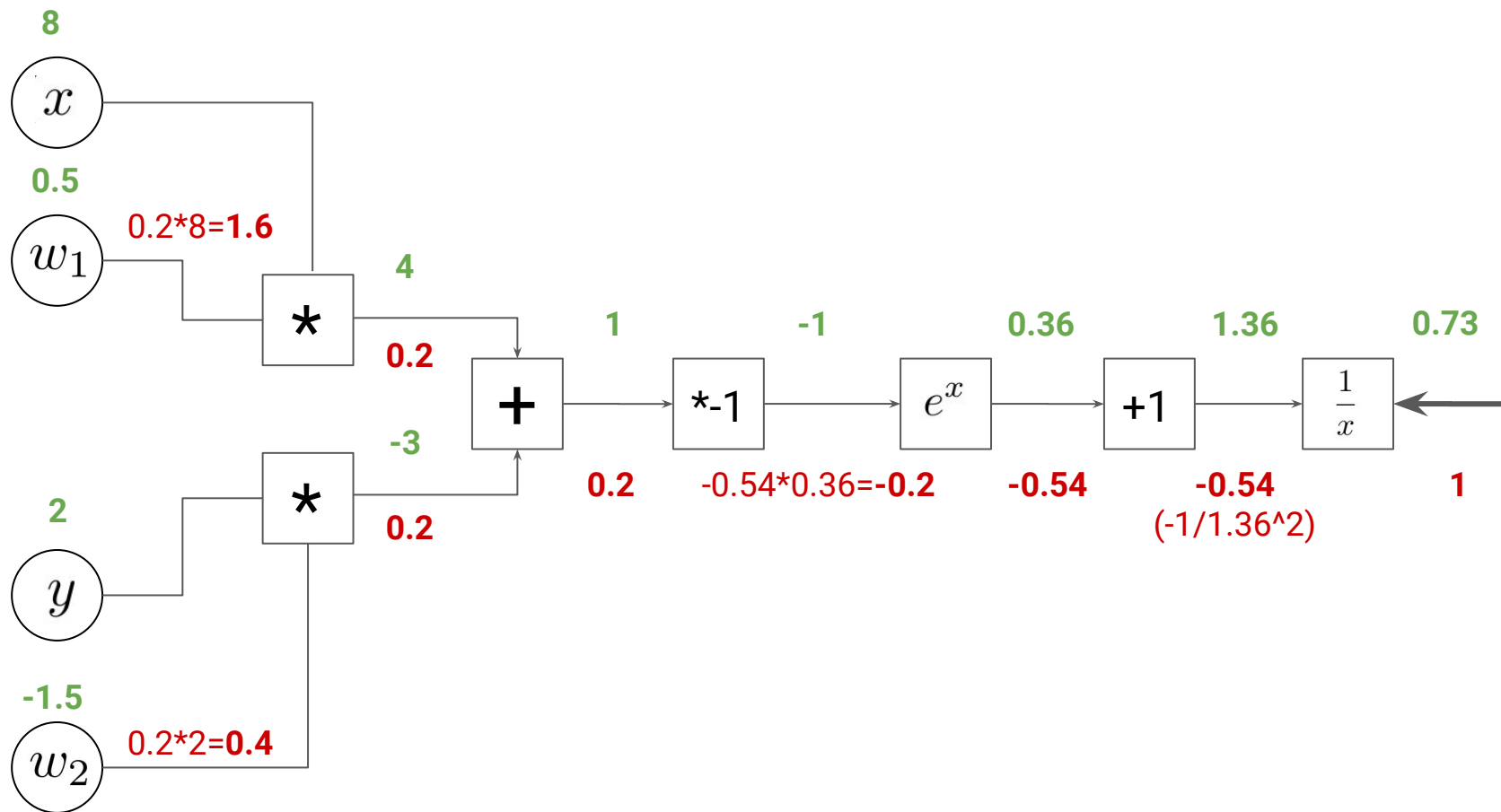
**Geoffrey Hinton**

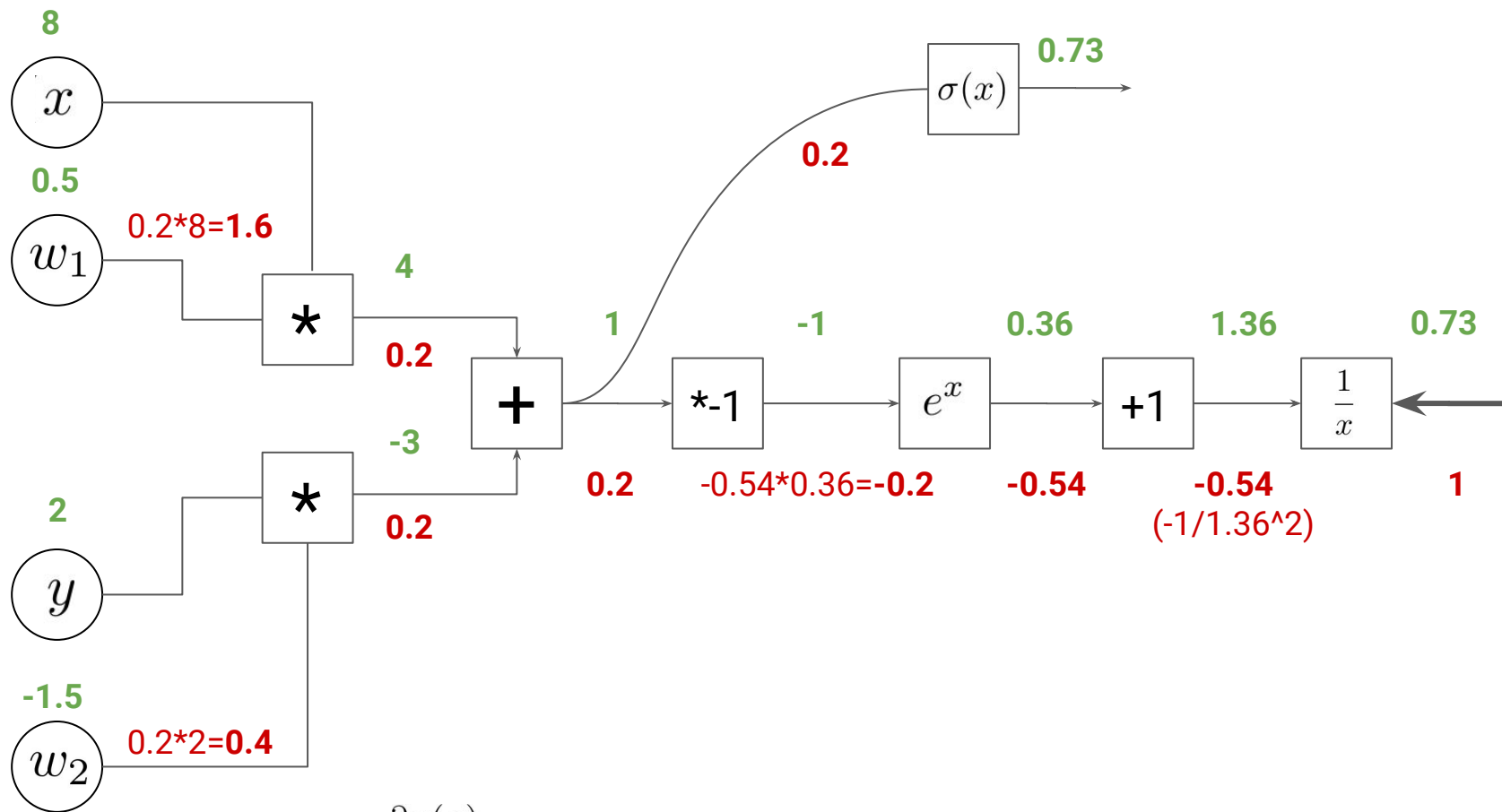
“Godfather of Deep Learning”











$$\frac{\partial \sigma(x)}{\partial x} = \sigma(x) (1 - \sigma(x))$$

# Стохастический градиентный спуск

# Градиентный спуск

1. Инициализировать параметры случайным образом (выбрать  $\theta_0$ )

2. Вычислить значение функционала  $\mathcal{L} = \frac{1}{n} \sum_{i=1}^n L(y_i, f(\mathbf{x}_i, \boldsymbol{\theta}))$

3. Вычислить  $\frac{\partial \mathcal{L}}{\partial \theta}(\theta_{i-1})$  при текущем значении параметра

4. Сделать шаг оптимизации:

$$\theta_i = \theta_{i-1} - \lambda \frac{\partial \mathcal{L}}{\partial \theta}(\theta_{i-1})$$

**а. После этого шага параметры изменились!**

5. Повторять 2-4 пока  $\mathcal{L}$  не перестанет меняться

# Градиентный спуск

1. Инициализировать параметры случайным образом (выбрать  $\theta_0$ )

2. Вычислить значение функционала  $\mathcal{L} = \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i, \theta))$

3. Вычислить  $\frac{\partial \mathcal{L}}{\partial \theta}(\theta_{i-1})$  при текущем значении параметра

4. Сделать шаг оптимизации:

$$\theta_i = \theta_{i-1} - \lambda \frac{\partial \mathcal{L}}{\partial \theta}(\theta_{i-1})$$

**а. После этого шага параметры изменились!**

5. Повторять 2-4 пока  $\mathcal{L}$  не перестанет меняться

# Градиентный спуск

1. Инициализировать параметры случайным образом (выбрать  $\theta_0$ )

2. Вычислить значение функционала 
$$\mathcal{L} = \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i, \theta))$$

3. Вычислить  $\frac{\partial \mathcal{L}}{\partial \theta}(\theta_{i-1})$  при текущем значении параметра

4. Сделать шаг оптимизации:

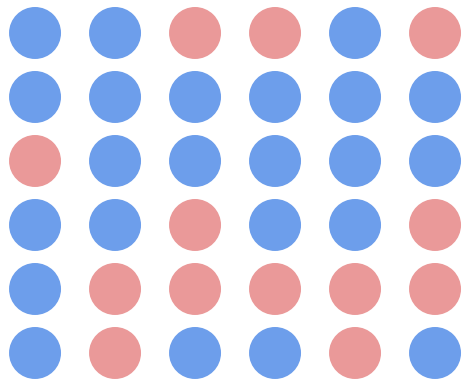
$$\theta_i = \theta_{i-1} - \lambda \frac{\partial \mathcal{L}}{\partial \theta}(\theta_{i-1})$$

**а. После этого шага параметры изменились!**

5. Повторять 2-4 пока  $\mathcal{L}$  не перестанет меняться

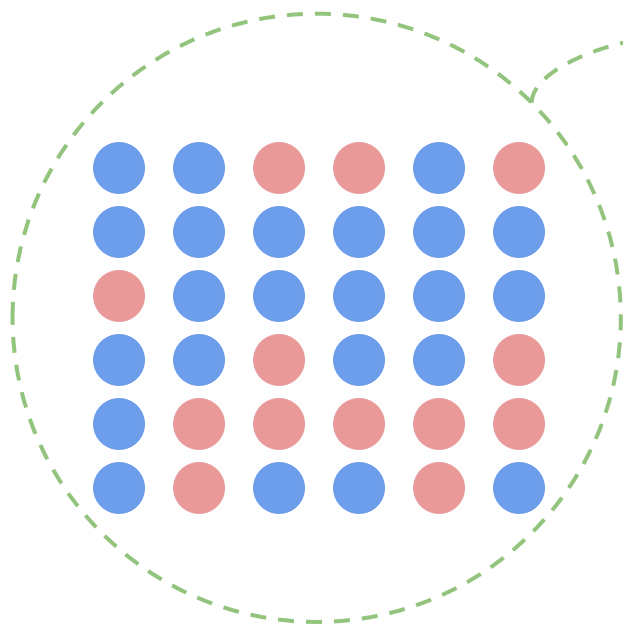
Нужно применить  $f$  ко всей выборке!  
Это может быть очень долго...

# Градиентный спуск





# Градиентный спуск

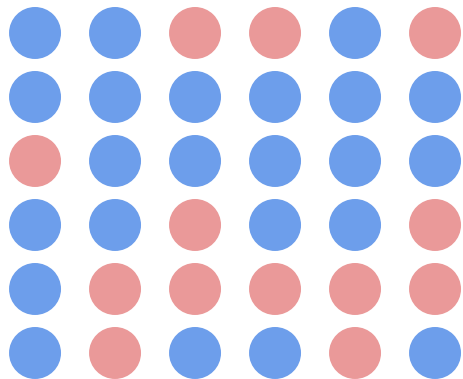


$$\mathcal{L} = \frac{1}{N} \sum_{i=1}^N L(y_i, \hat{y}_i)$$

+ Шаг оптимизации

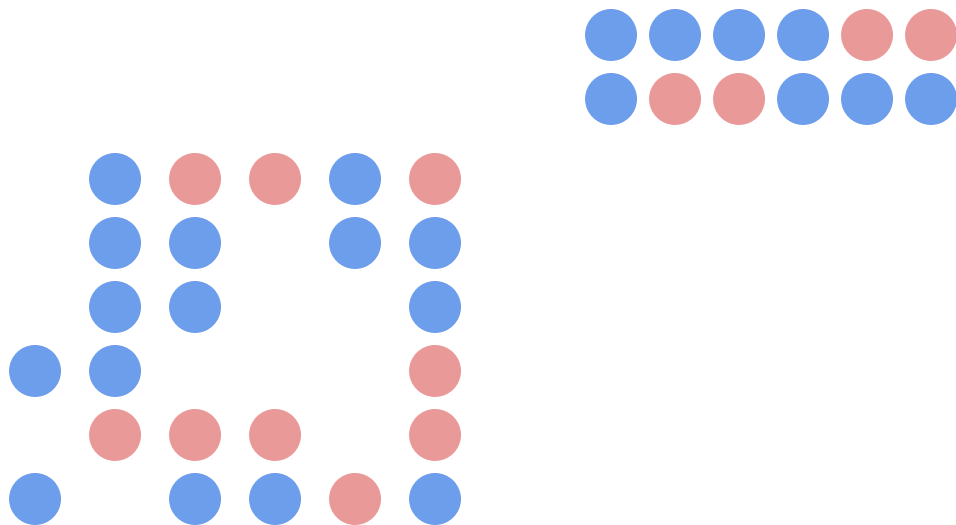
# Стохастический градиентный спуск

Stochastic Gradient Descent (SGD)



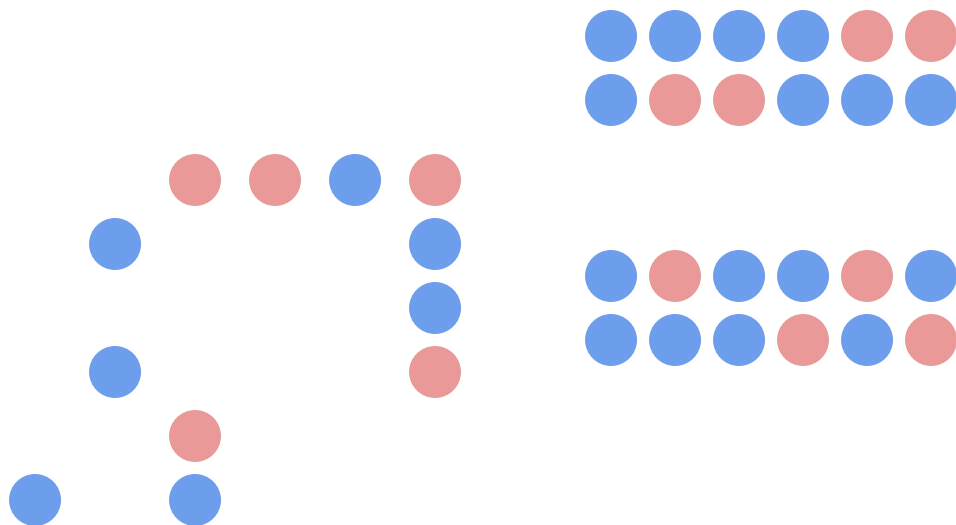
# Стохастический градиентный спуск

Stochastic Gradient Descent (SGD)



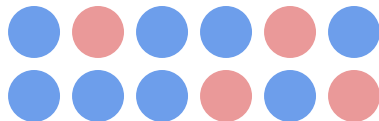
# Стохастический градиентный спуск

Stochastic Gradient Descent (SGD)



# Стохастический градиентный спуск

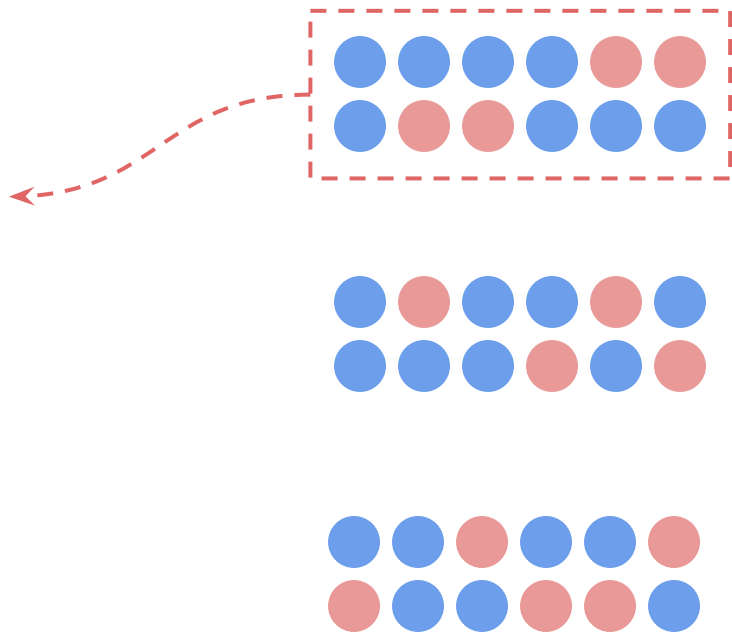
Stochastic Gradient Descent (SGD)



# Стохастический градиентный спуск

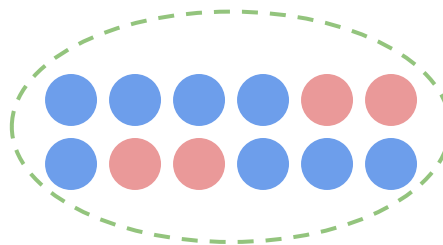
Stochastic Gradient Descent (SGD)

Batch, mini-batch  
Batch size = 12



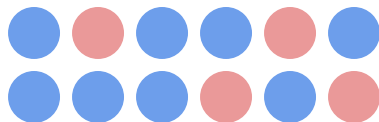
# Стохастический градиентный спуск

Stochastic Gradient Descent (SGD)



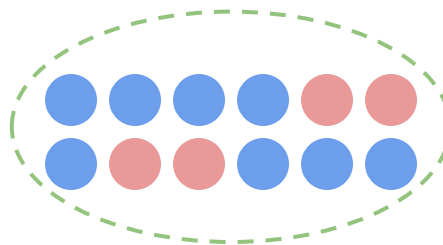
A diagram showing a batch of 12 data points arranged in two rows of six. The top row contains four blue circles followed by two red circles. The bottom row contains one blue circle, two red circles, and three blue circles. A dashed green oval encloses the entire batch. A dashed green arrow points from the oval to the loss function equation.

$$\mathcal{L} = \frac{1}{n_1} \sum_{i \in [\text{batch 1 ids}]} L(y_i, \hat{y}_i)$$



# Стохастический градиентный спуск

Stochastic Gradient Descent (SGD)



A diagram showing a batch of 10 data points arranged in two rows of five. The top row contains four blue circles followed by two red circles. The bottom row contains one blue circle, two red circles, and three blue circles. A dashed green oval encloses the entire batch. A dashed green arrow points from the right side of the oval to the loss function equation.

$$\mathcal{L} = \frac{1}{n_1} \sum_{i \in [\text{batch 1 ids}]} L(y_i, \hat{y}_i)$$

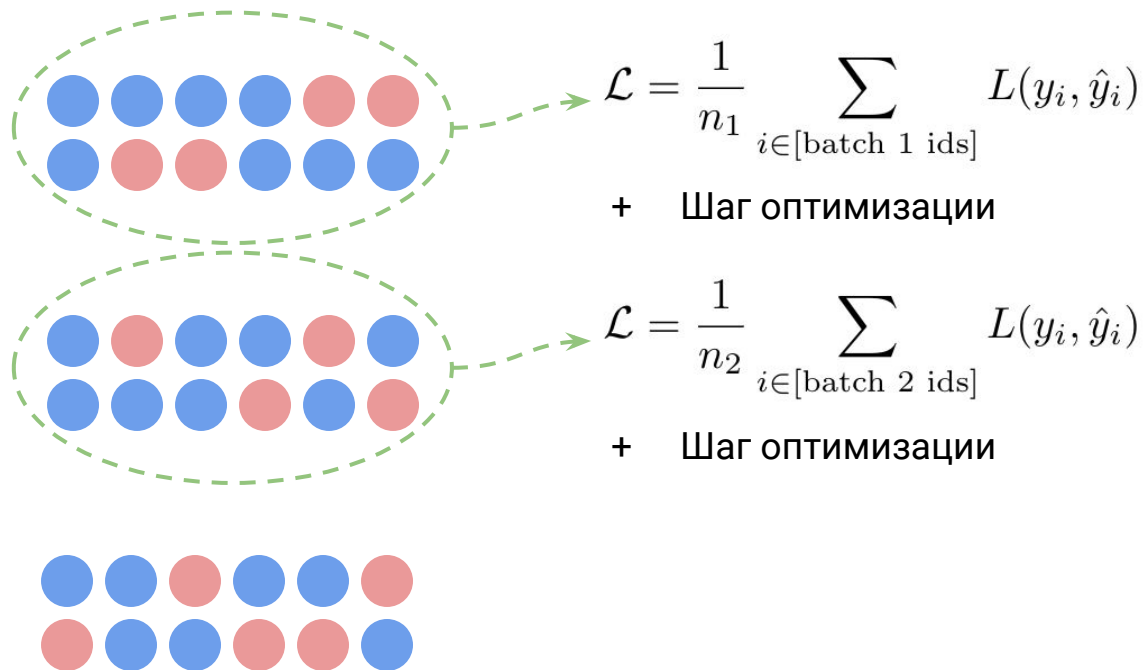
+ Шаг оптимизации





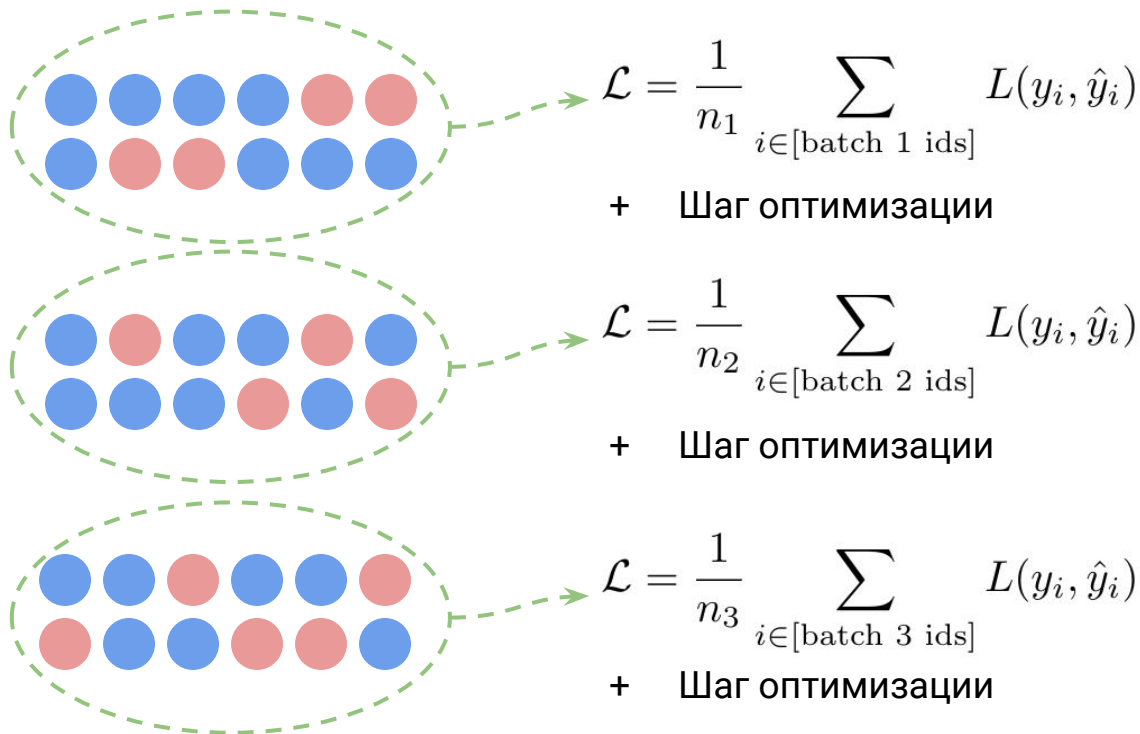
# Стохастический градиентный спуск

Stochastic Gradient Descent (SGD)



# Стохастический градиентный спуск

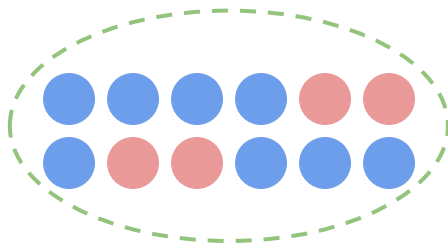
Stochastic Gradient Descent (SGD)



# Стохастический градиентный спуск

Stochastic Gradient Descent (SGD)

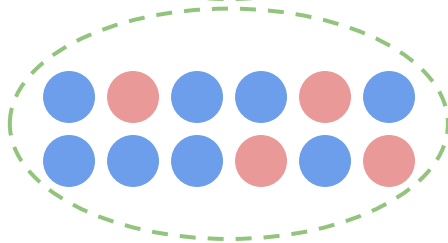
Итерация №1



$$\mathcal{L} = \frac{1}{n_1} \sum_{i \in [\text{batch 1 ids}]} L(y_i, \hat{y}_i)$$

+ Шаг оптимизации

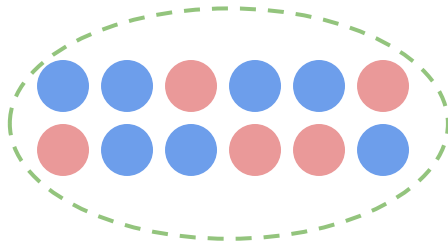
Итерация №2



$$\mathcal{L} = \frac{1}{n_2} \sum_{i \in [\text{batch 2 ids}]} L(y_i, \hat{y}_i)$$

+ Шаг оптимизации

Итерация №3

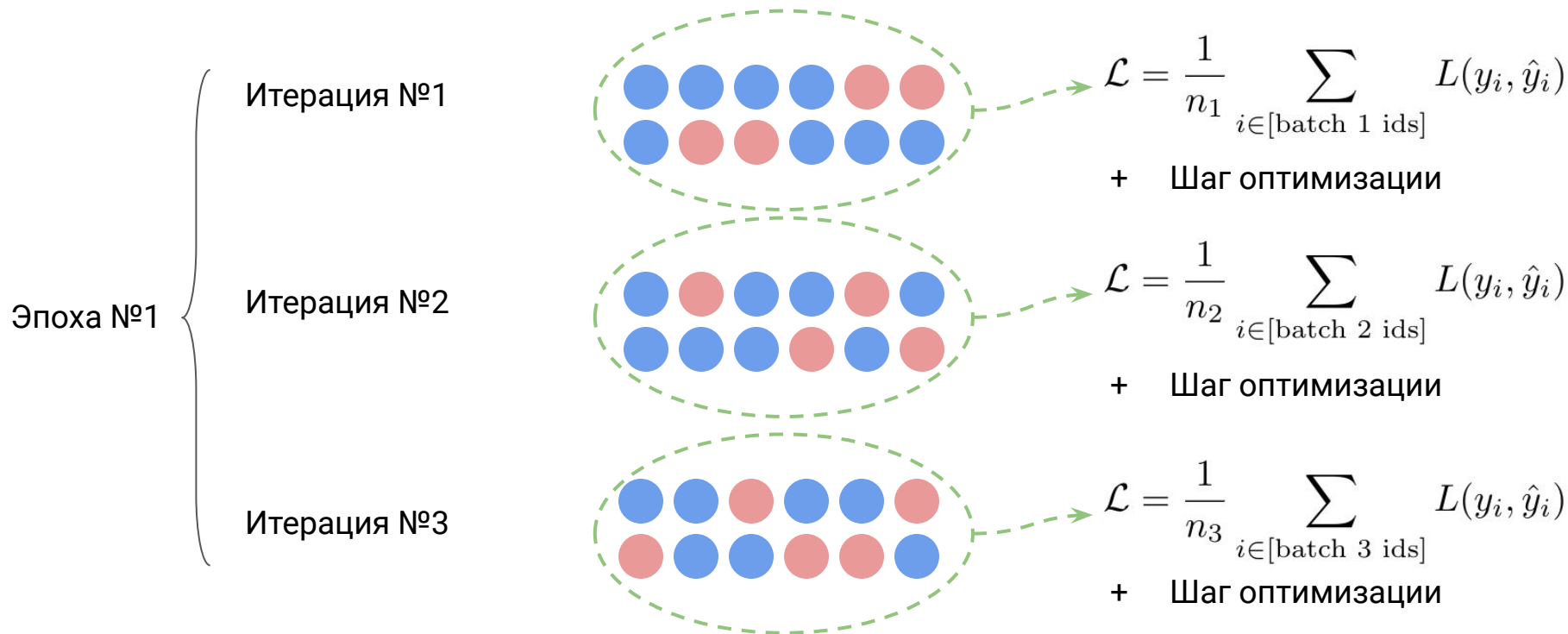


$$\mathcal{L} = \frac{1}{n_3} \sum_{i \in [\text{batch 3 ids}]} L(y_i, \hat{y}_i)$$

+ Шаг оптимизации

# Стохастический градиентный спуск

Stochastic Gradient Descent (SGD)

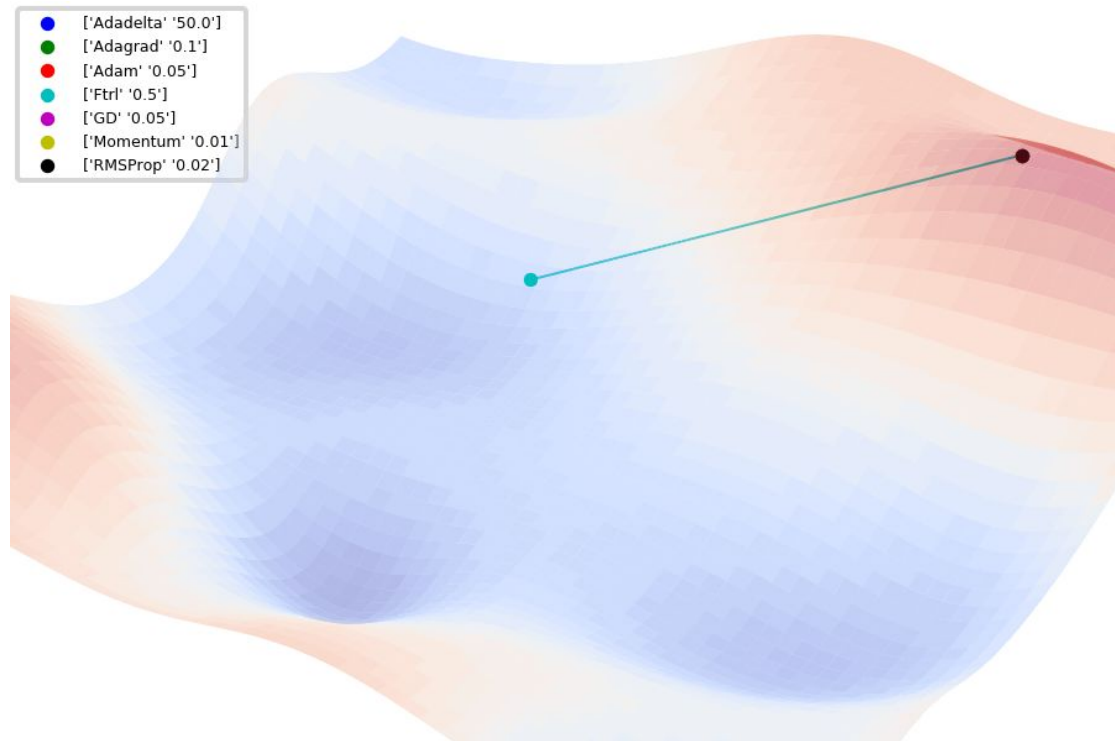


# Достоинства и недостатки SGD

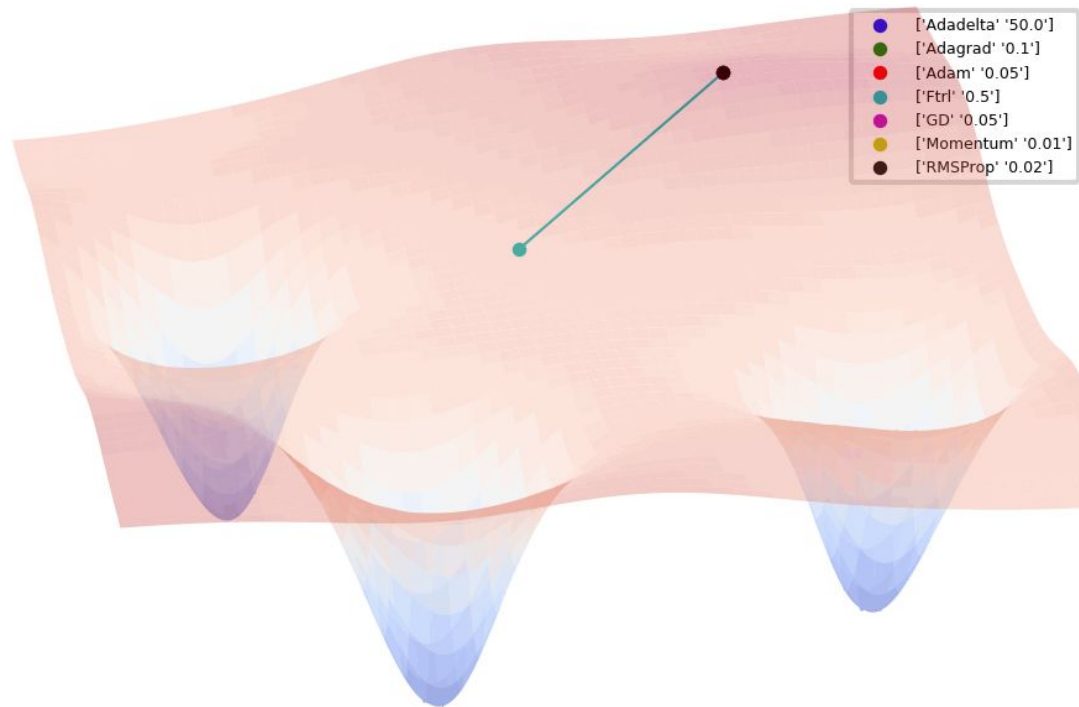
- Достоинства:
  - Позволяет делать шаги оптимизации чаще
  - Легко реализуется
- Недостатки:
  - Застревание в локальных минимумах
  - Медленная сходимость

На практике чаще используются модификации SGD (например, Adam)

# Модификации градиентного спуска



# Модификации градиентного спуска



# ADAM (упрощенно)

GD update:

$$\theta_t \leftarrow \theta_{t-1} - \alpha \nabla_{\theta} f_t(\theta_{t-1})$$



# ADAM (упрощенно)

GD update:  $\theta_t \leftarrow \theta_{t-1} - \alpha \nabla_{\theta} f_t(\theta_{t-1})$

ADAM update:  $\theta_t \leftarrow \theta_{t-1} - \alpha \cdot \hat{m}_t / (\sqrt{\hat{v}_t} + \epsilon)$

# ADAM (упрощенно)

GD update:  $\theta_t \leftarrow \theta_{t-1} - \alpha \nabla_{\theta} f_t(\theta_{t-1})$

ADAM update:  $\theta_t \leftarrow \theta_{t-1} - \alpha \cdot \hat{m}_t / (\sqrt{\hat{v}_t} + \epsilon)$

Экспоненциальное скользящее  
среднее градиента  
(“идем примерно туда куда шли”)

Экспоненциальное скользящее  
среднее квадрата градиента  
(“если сильно болтает -- делаем  
шаги меньше”)

# Итог

- Классический градиентный спуск требует расчета ошибки по всей выборке данных для того чтобы сделать шаг оптимизации
- Стохастический градиентный спуск работает с “кусочками” данных, делая обновления чаще
- Тем не менее он склонен к застреванию в локальных минимумах (как и классический вариант)
- Есть множество модификаций алгоритма, с одним из самых популярных из которых мы познакомились на лекции

# Функции активации

# Нейрон

$$\sigma(wx + b) = p_+$$

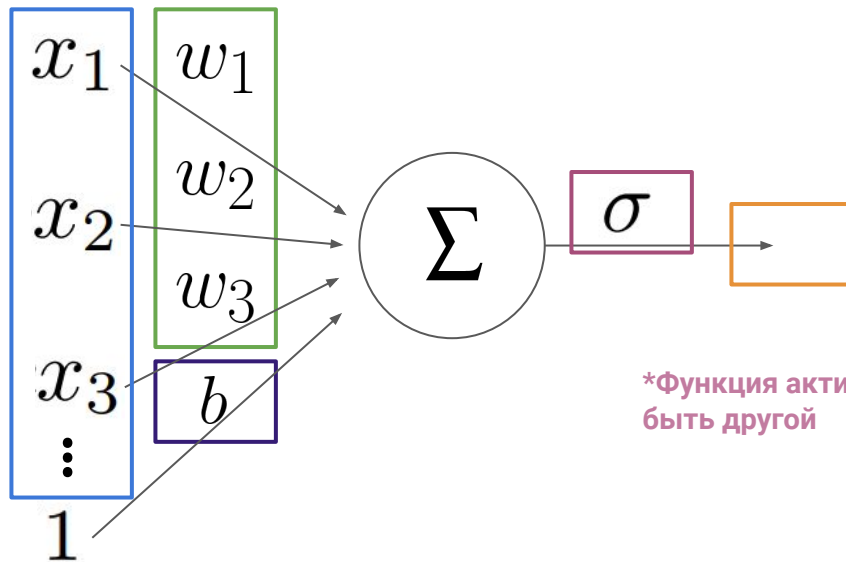
Вход нейрона  
(признаки, inputs)

Коэффициенты (веса,  
weights)

Смещение (bias)

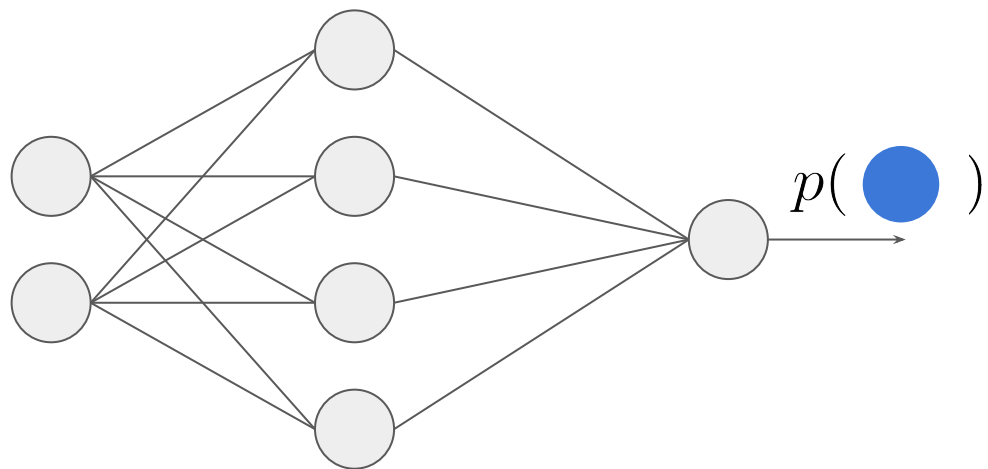
Функция активации

Выход нейрона

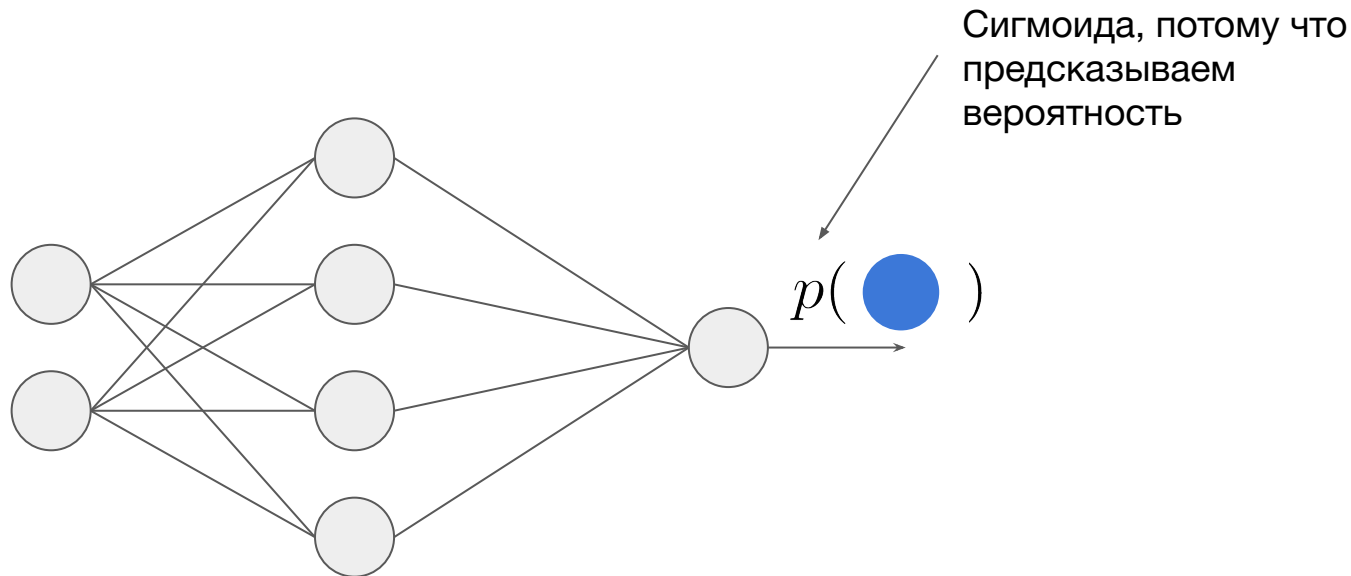


\*Функция активации может  
быть другой

# Функции активации

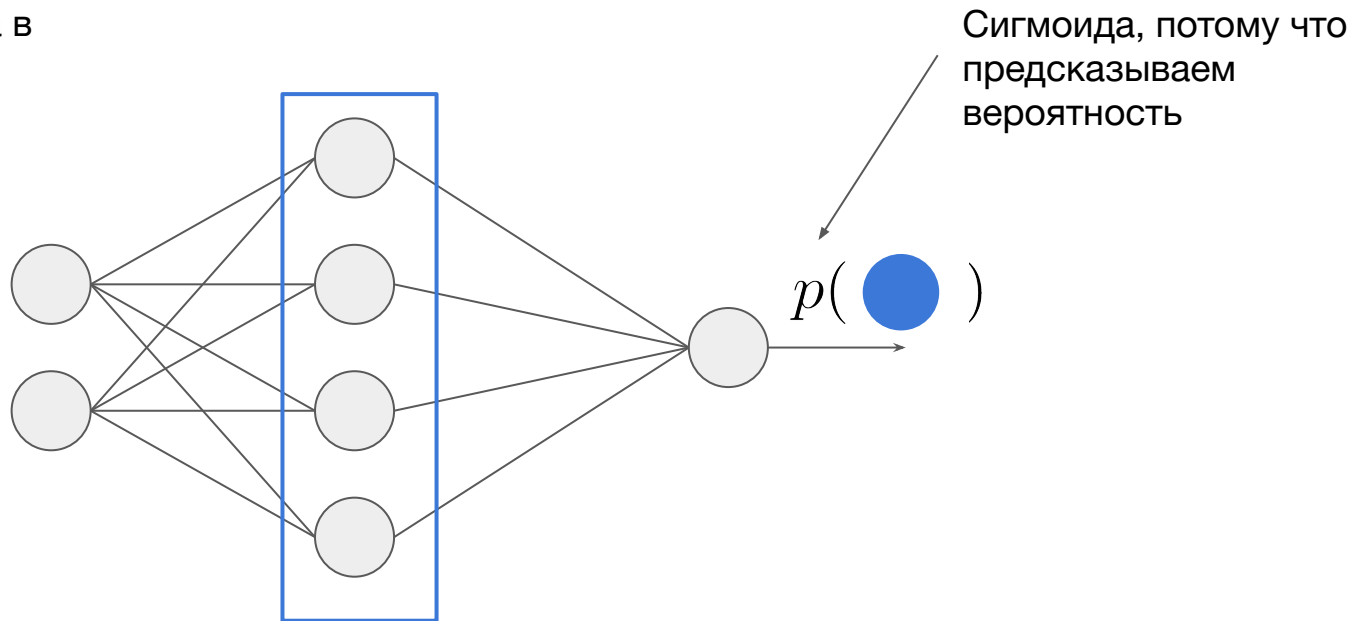


# Функции активации



# Функции активации

Обязательна ли сигмоида в промежуточных слоях?

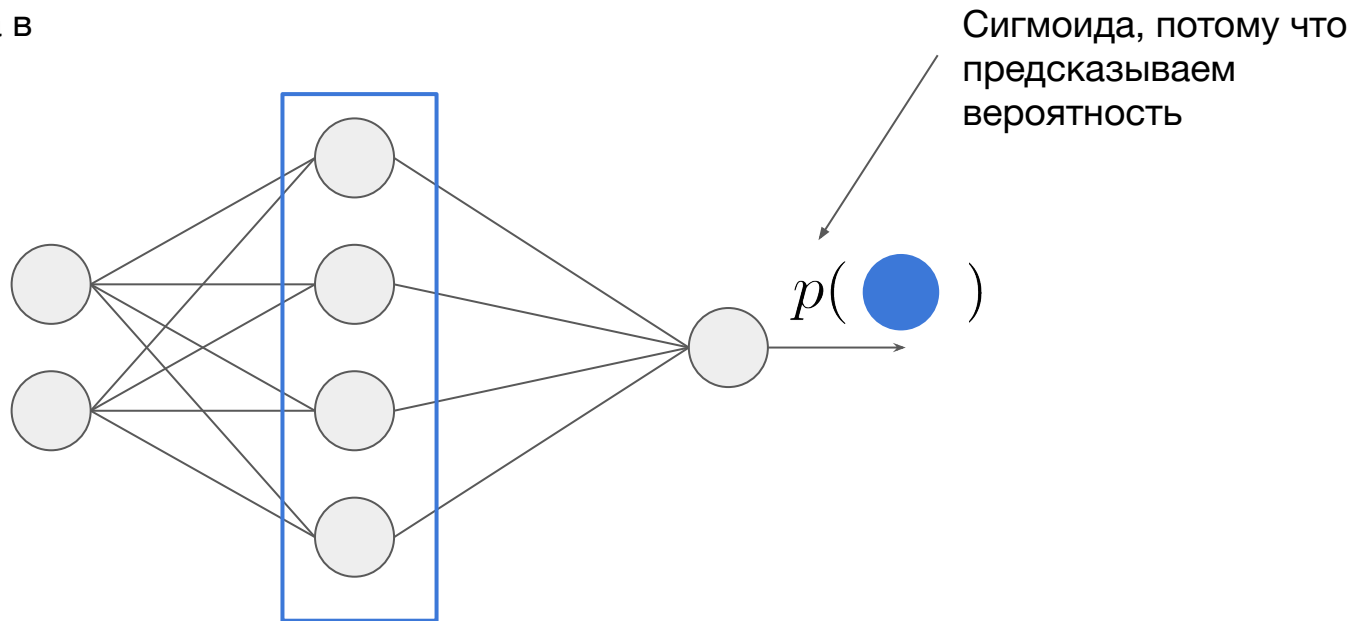




# Функции активации

Обязательна ли сигмоида в промежуточных слоях?

- **Нет**

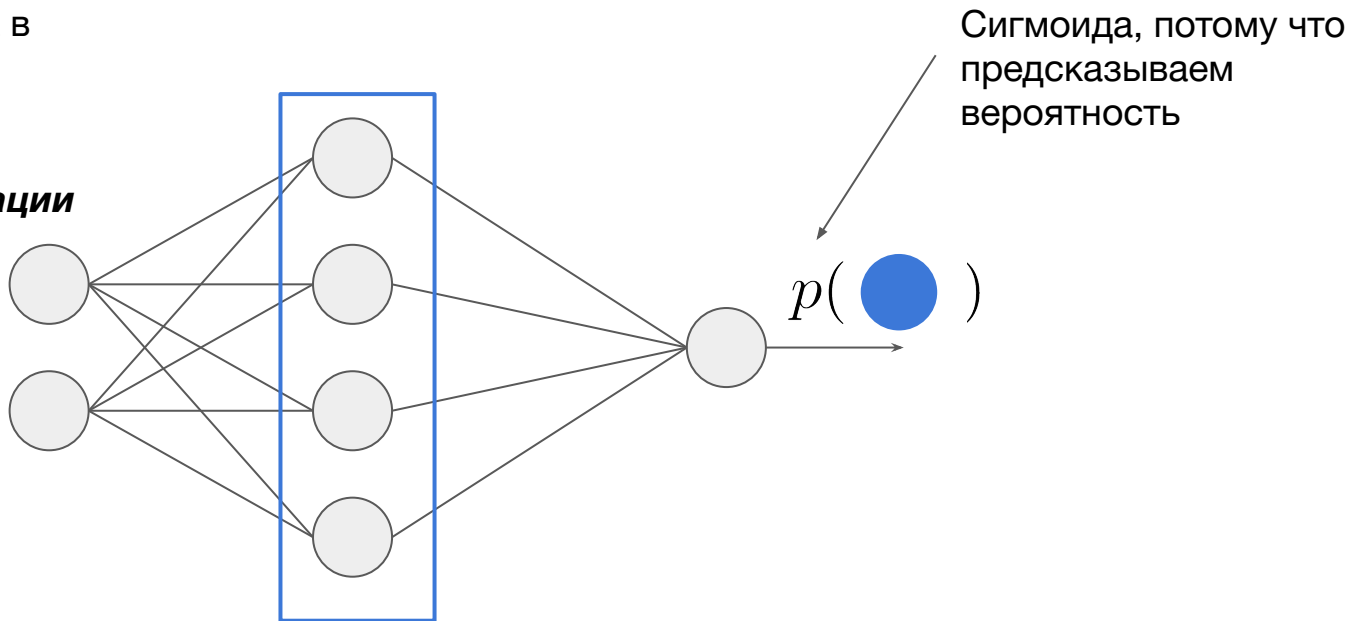


# Функции активации

Обязательна ли сигмоида в промежуточных слоях?

- Нет

*Можно ли убрать активации вообще?*



# Функции активации

$$xw_1 + b_1$$

Линейная функция (один нейрон без активации, один признак на входе)

# Функции активации

$$\sigma((xw_1 + b_1)w_2 + b_2)$$

Применили второй нейрон и сигмоиду, чтобы получить вероятность

# Функции активации

$$\sigma((xw_1 + b_1)w_2 + b_2)$$

Применили второй нейрон и сигмоиду, чтобы получить вероятность  
Сделали мы функцию сложнее?

# Функции активации

$$\sigma((xw_1 + b_1)w_2 + b_2)$$

$$xw_1w_2 + b_1w_2 + b_2$$

Применили второй нейрон и сигмоиду, чтобы получить вероятность  
Сделали мы функцию сложнее?

# Функции активации

$$\sigma((xw_1 + b_1)w_2 + b_2)$$

$$w = w_1w_2, b = b_1w_2 + b_2 \quad x\boxed{w_1w_2} + \boxed{b_1w_2 + b_2}$$

Применили второй нейрон и сигмоиду, чтобы получить вероятность  
Сделали мы функцию сложнее?

# Функции активации

$$\sigma((xw_1 + b_1)w_2 + b_2)$$

$$w = w_1w_2, b = b_1w_2 + b_2 \quad x\boxed{w_1w_2} + \boxed{b_1w_2 + b_2} \longrightarrow xw + b$$

Применили второй нейрон и сигмоиду, чтобы получить вероятность  
Сделали мы функцию сложнее?



# Функции активации

$$\sigma((xw_1 + b_1)w_2 + b_2)$$

$$w = w_1w_2, b = b_1w_2 + b_2 \quad x \boxed{w_1w_2} + \boxed{b_1w_2 + b_2} \longrightarrow xw + b$$

Применили второй нейрон и сигмоиду, чтобы получить вероятность

Сделали мы функцию сложнее?

Нет! Линейная функция от линейной функции -- линейная функция!

# Функции активации

$$\sigma((xw_1 + b_1)w_2 + b_2)$$

$$w = w_1w_2, b = b_1w_2 + b_2 \quad x\boxed{w_1w_2} + \boxed{b_1w_2 + b_2} \longrightarrow xw + b$$

Применили второй нейрон и сигмоиду, чтобы получить вероятность

Сделали мы функцию сложнее?

Нет! Линейная функция от линейной функции -- линейная функция!

**Чтобы придумывать более сложные признаки нам нужна нелинейность**

# Функции активации

$$\sigma((xw_1 + b_1)w_2 + b_2)$$

$$w = w_1w_2, b = b_1w_2 + b_2 \quad x\boxed{w_1w_2} + \boxed{b_1w_2 + b_2} \longrightarrow xw + b$$

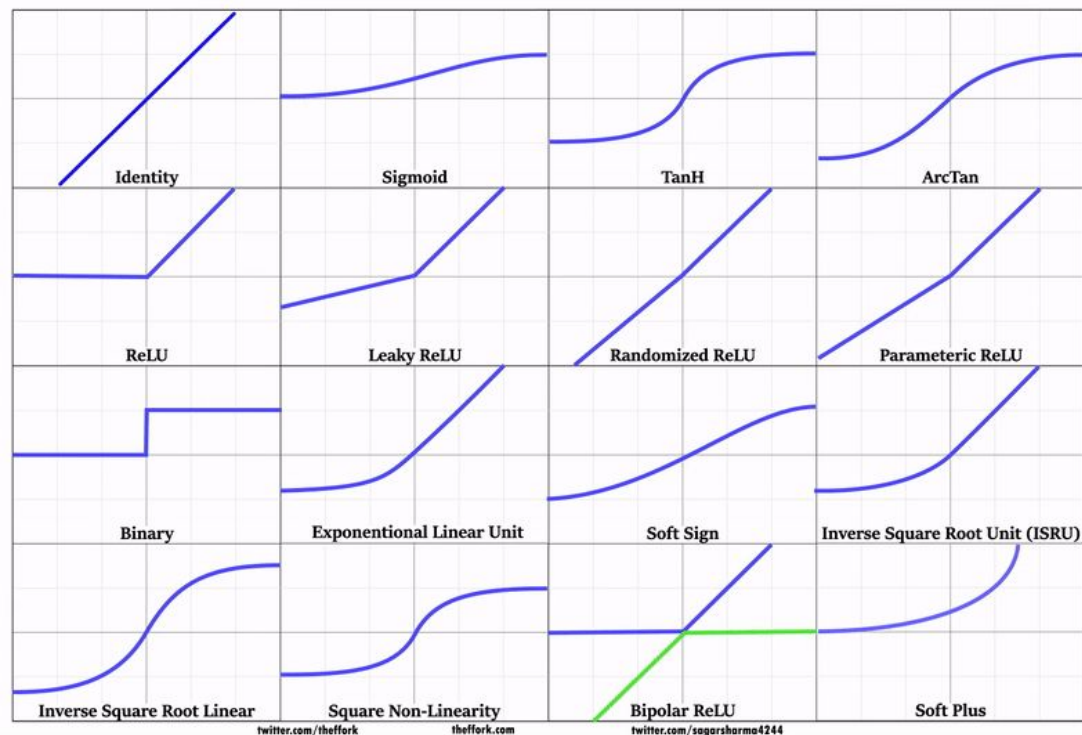
Применили второй нейрон и сигмоиду, чтобы получить вероятность

Сделали мы функцию сложнее?

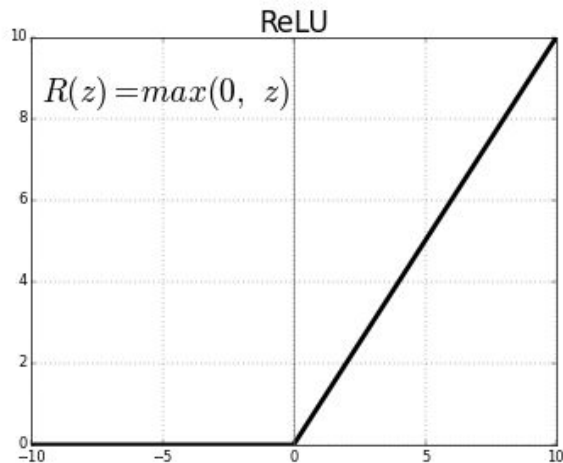
Нет! Линейная функция от линейной функции -- линейная функция!

**Чтобы придумывать более сложные признаки нам нужна нелинейность  
Но необязательно сигмоида)**

# Функции активации



# ReLU



- Основная функции активации на сегодняшний день
- Она нелинейна
- Чем она лучше сигмоиды?
  - Короткий ответ: у нее лучше производная
  - Она быстрее считается
- Но сигмоиду все равно используем к концу, чтобы получить вероятность!

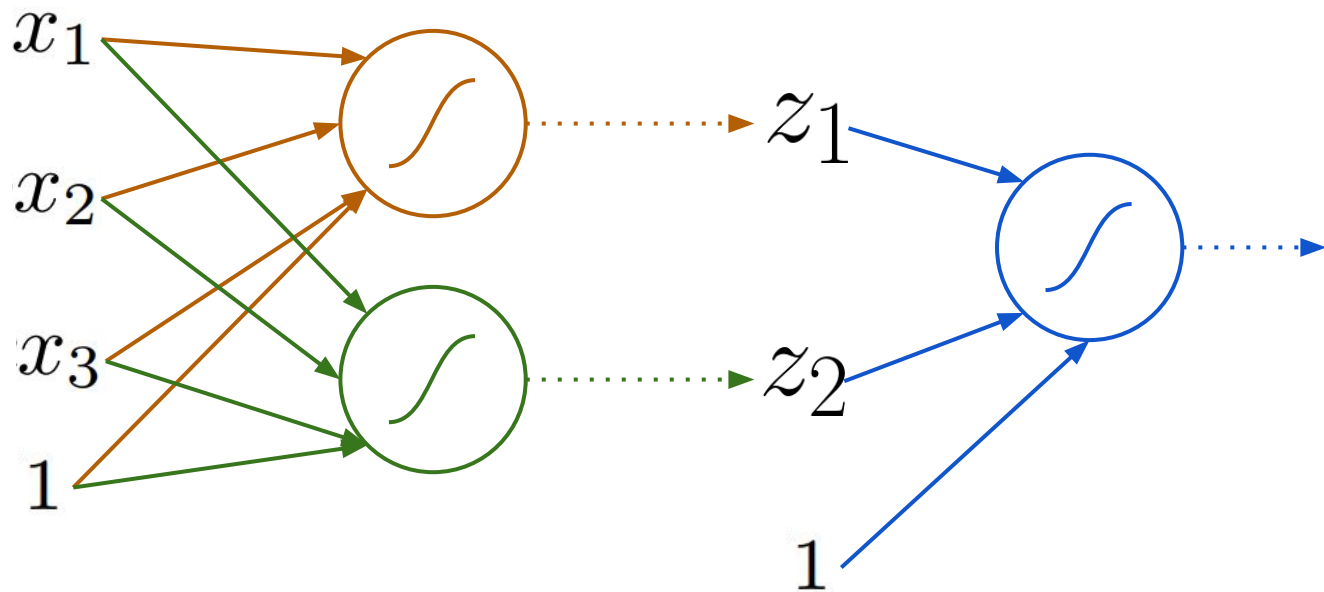
Обучение нейронных сетей

# Нейронные сети: от схем к формулам

# Как задать полносвязный слой в коде

- На практике прошлого блока мы реализовали нейронную сеть из трех нейронов
- Каждый мы определяли отдельно
- Как быть, если нейронов много? (а их всегда много)
- Для начала поймем как упростить запись

# Нейронная сеть



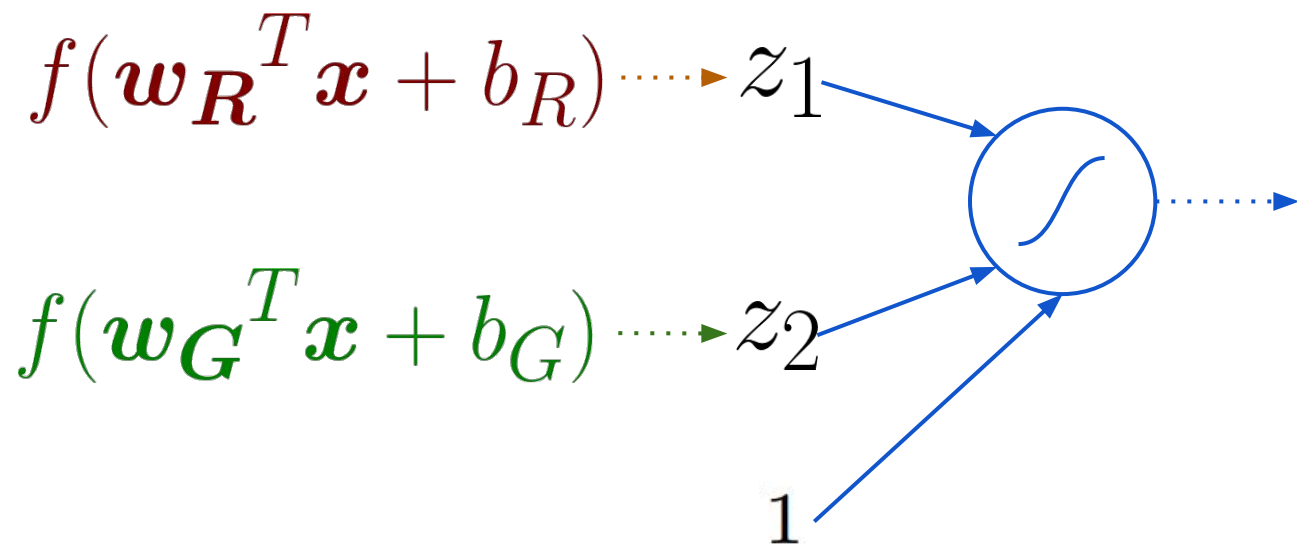


# Нейронная сеть

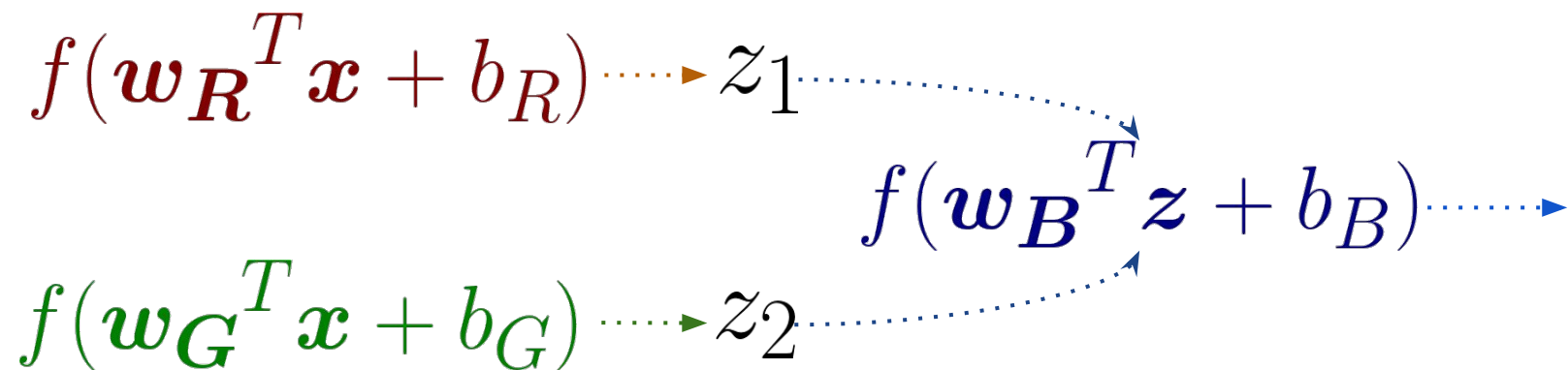
$w$  - вектор размера  $m$

$x$  - вектор размера  $m$

$b$  - число



# Нейронная сеть



$$f(\boldsymbol{w}_R^T \boldsymbol{x} + b_R)$$

$$f(\boldsymbol{w}_G^T \boldsymbol{x} + b_G)$$

$$f(\boldsymbol{w}_R^T \boldsymbol{x} + b_R)$$

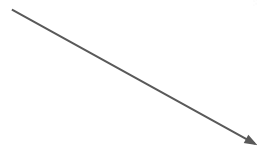
$$(w_R^1 \quad w_R^2 \quad \dots \quad w_R^m) \begin{pmatrix} x_1 \\ x_2 \\ \dots \\ x_m \end{pmatrix} + b_R$$

$$f(\boldsymbol{w}_G^T \boldsymbol{x} + b_G)$$

$$(w_G^1 \quad w_G^2 \quad \dots \quad w_G^m) \begin{pmatrix} x_1 \\ x_2 \\ \dots \\ x_m \end{pmatrix} + b_G$$

$$f(\mathbf{w}_R^T \mathbf{x} + b_R)$$

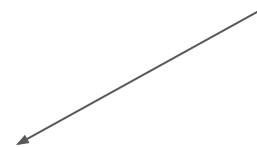
$$(w_R^1 \ w_R^2 \ \dots \ w_R^m) \begin{pmatrix} x_1 \\ x_2 \\ \dots \\ x_m \end{pmatrix} + b_R$$



$$\begin{pmatrix} w_R^1 & w_R^2 & \dots & w_R^m \\ w_G^1 & w_G^2 & \dots & w_G^m \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \dots \\ x_m \end{pmatrix} + \begin{pmatrix} b_R \\ b_G \end{pmatrix}$$

$$f(\mathbf{w}_G^T \mathbf{x} + b_G)$$

$$(w_G^1 \ w_G^2 \ \dots \ w_G^m) \begin{pmatrix} x_1 \\ x_2 \\ \dots \\ x_m \end{pmatrix} + b_G$$




$$f(\mathbf{w}_R^T \mathbf{x} + b_R)$$

$$f(\mathbf{w}_G^T \mathbf{x} + b_G)$$

$$(w_R^1 \quad w_R^2 \quad \dots \quad w_R^m) \begin{pmatrix} x_1 \\ x_2 \\ \dots \\ x_m \end{pmatrix} + b_R$$


$$(w_G^1 \quad w_G^2 \quad \dots \quad w_G^m) \begin{pmatrix} x_1 \\ x_2 \\ \dots \\ x_m \end{pmatrix} + b_G$$



The diagram shows two arrows pointing from the equations above to a single equation below. The left arrow points from the red equation, and the right arrow points from the green equation. The resulting equation below has a matrix where the top row is red and the bottom row is green, followed by a vector  $\mathbf{x}$  and a vector  $\mathbf{b}$  where the top element is red and the bottom is green.

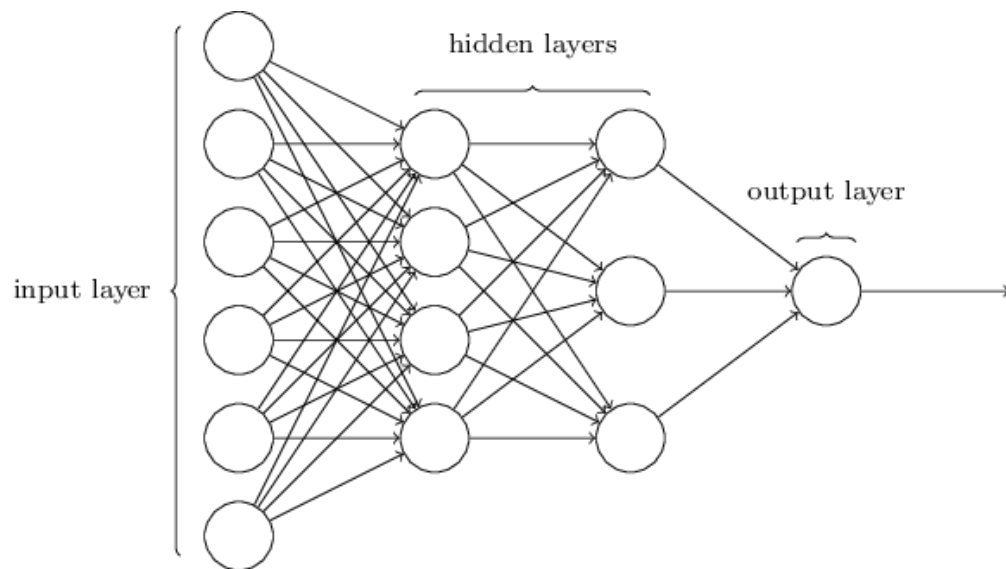
$$\begin{pmatrix} w_R^1 & w_R^2 & \dots & w_R^m \\ w_G^1 & w_G^2 & \dots & w_G^m \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ \dots \\ x_m \end{pmatrix} + \begin{pmatrix} b_R \\ b_G \end{pmatrix}$$

$$W\mathbf{x} + \mathbf{b}$$

$$Wx + b$$


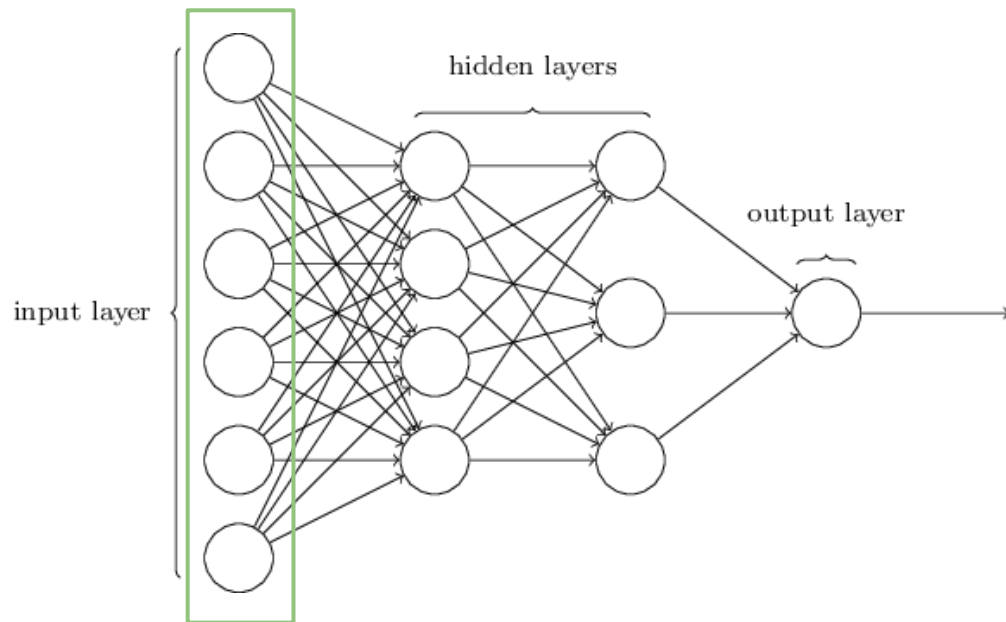
- Матрица весов: (2, m) в нашем примере
  - (k, m) – в общем случае.
  - m – количество нейронов на входе в слой, k – на выходе
- Вектор смещения: длины 2
  - k – в общем случае

# Нейронная сеть

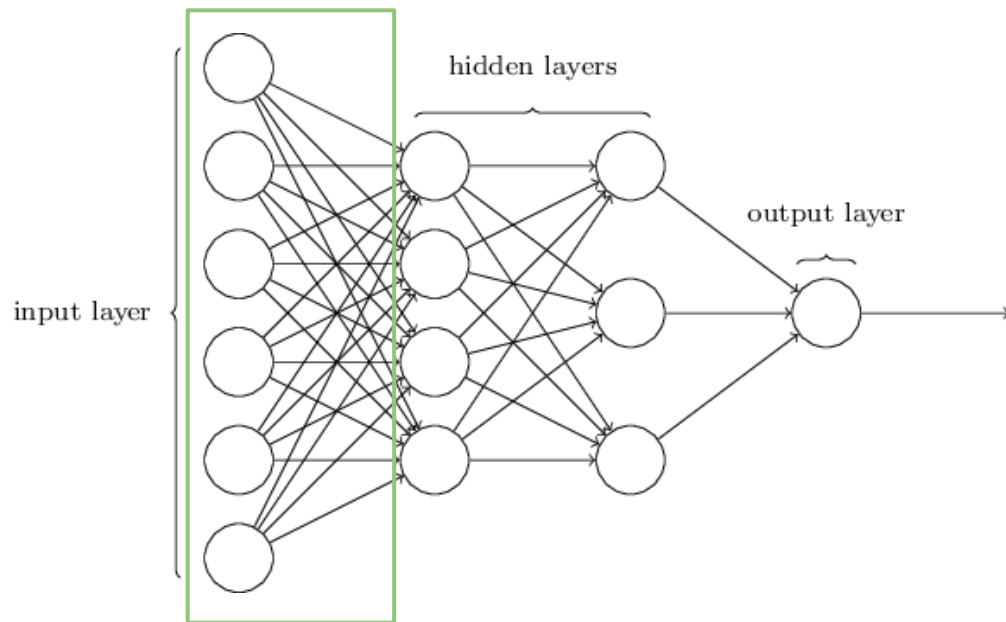




# Нейронная сеть

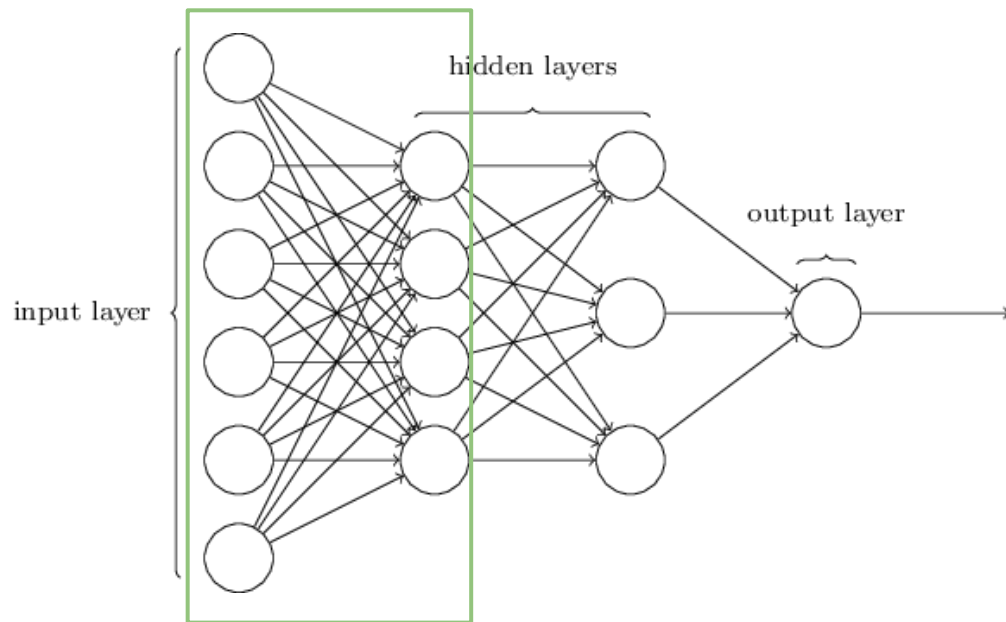
 $x$

# Нейронная сеть



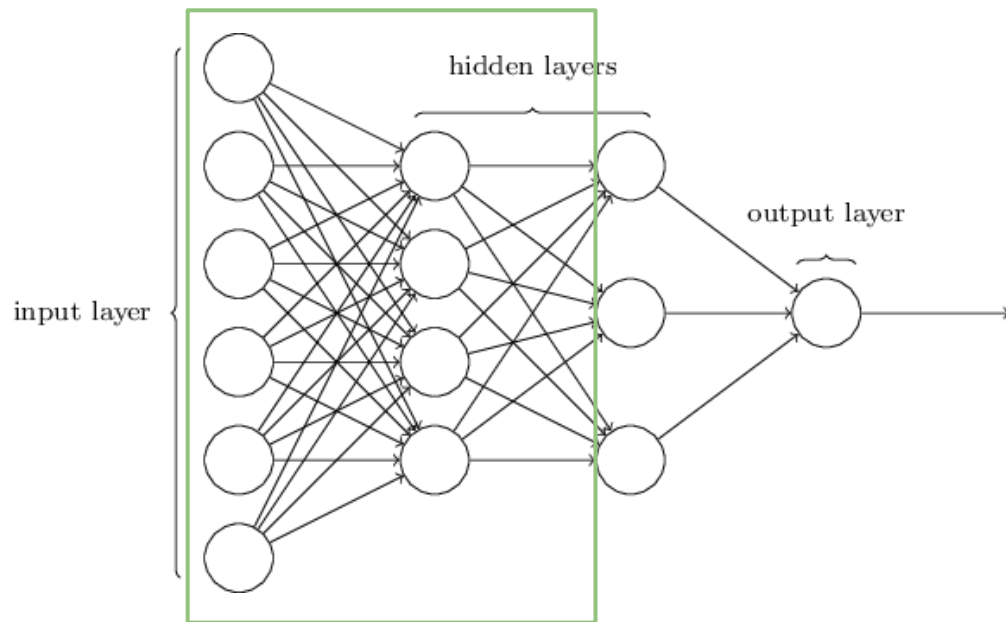
$$W_1x + b_1$$

# Нейронная сеть



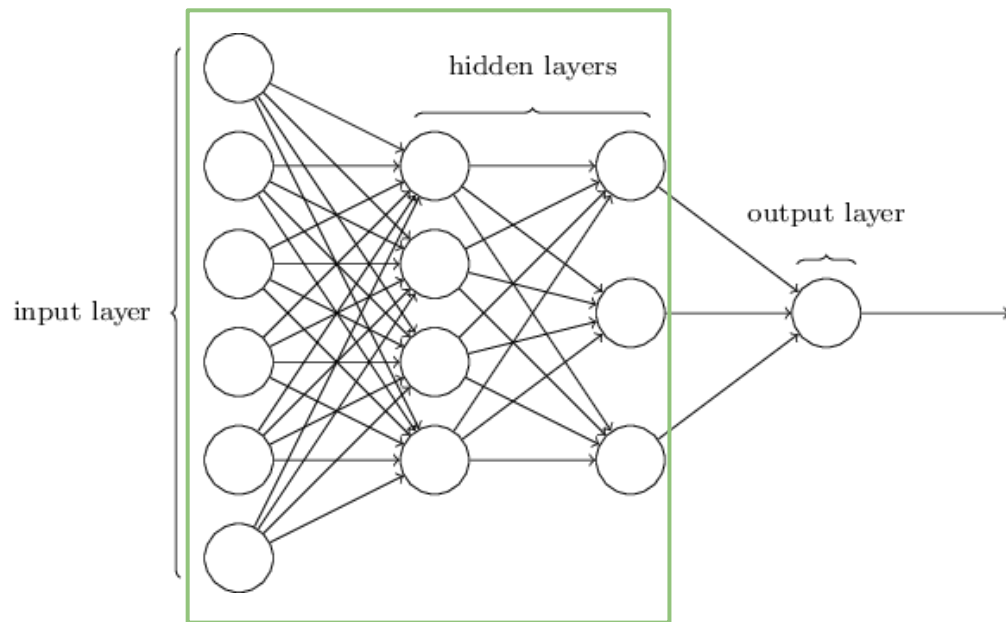
$$f(W_1x + b_1)$$

# Нейронная сеть



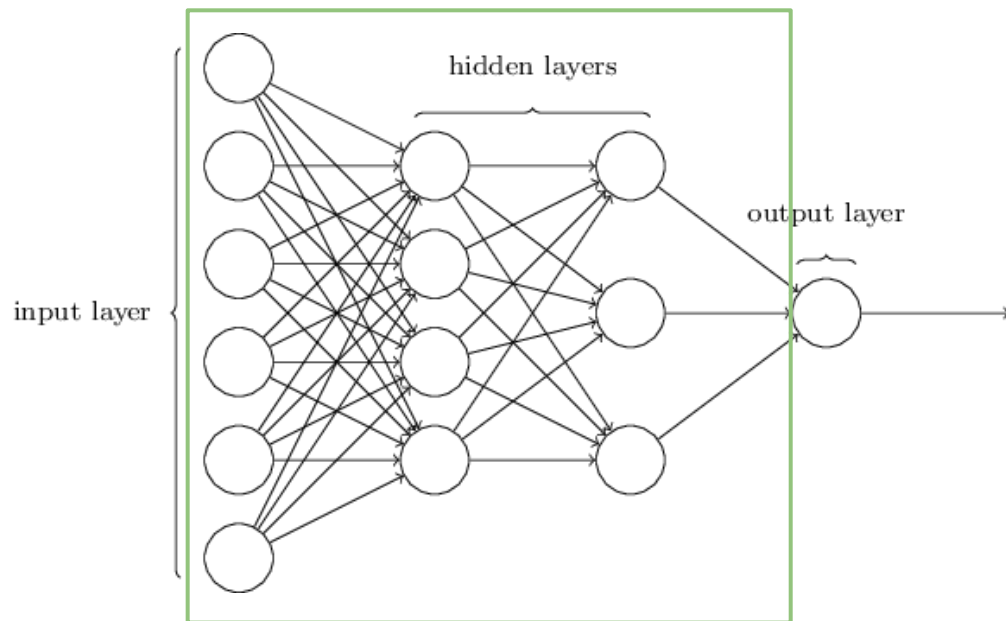
$$W_2 f(W_1 x + b_1) + b_2$$

# Нейронная сеть



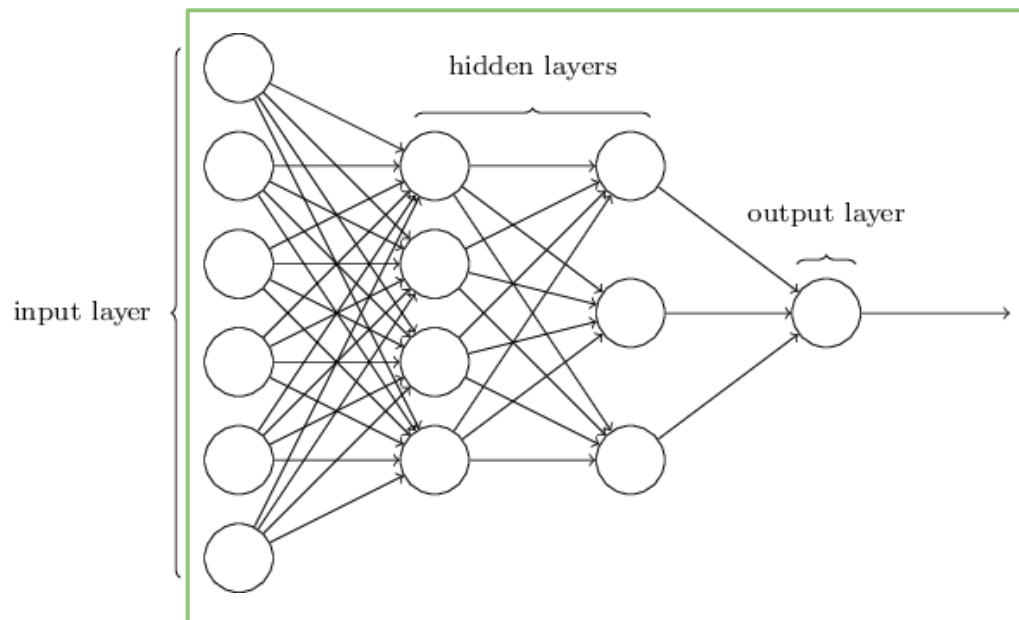
$$f(W_2 f(W_1 x + b_1) + b_2)$$

# Нейронная сеть



$$W_3(f(W_2f(W_1x + b_1) + b_2) + b_3$$

# Нейронная сеть



$$f(W_3(f(W_2f(W_1x + b_1) + b_2) + b_3))$$

# Полносвязный слой на numpy

```
m = 10
k = 2

x = np.ones(m)
W = np.ones((k, m))
b = np.ones(k)
print(x.shape, W.shape, b.shape)
```

```
(10,) (2, 10) (2,)
```



# Полносвязный слой на numpy

```
m = 10
k = 2

x = np.ones(m)
W = np.ones((k, m))
b = np.ones(k)
print(x.shape, W.shape, b.shape)
```

```
(10,) (2, 10) (2,)
```

```
np.dot(W, x) + b
```

```
array([11., 11.])
```

# Как применить к нескольким объектам?

```
n = 3
m = 10
k = 2

x = np.ones((n, m))
W = np.ones((m, k))
b = np.ones(k)
print(x.shape, W.shape, b.shape)
```

Batch dimension



```
(3, 10) (10, 2) (2,)
```

# Как применить к нескольким объектам?

```
n = 3
m = 10
k = 2

x = np.ones((n, m))
W = np.ones((m, k))
b = np.ones(k)
print(x.shape, W.shape, b.shape)

(3, 10) (10, 2) (2,)
```

```
np.matmul(x, W) + b

array([[11., 11.],
       [11., 11.],
       [11., 11.]])
```

# Итог

- Применение полносвязного слоя эквивалентно матричному умножению входа на матрицу весов и добавление вектора смещения
- Соответственно применимы все матричные операции и векторизация в numpy
- Например посчитать выход слоя для нескольких объектов можно сразу

Обучение нейронных сетей

# Мультиклассовая классификация

# Домашнее задание



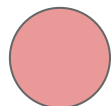
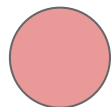
- Классификация рукописных цифр
- Датасет -- MNIST
- 60 000 изображений в трейне
- **10 классов**

# Домашнее задание



- Классификация рукописных цифр
- Датасет -- MNIST
- 60 000 изображений в трейне
- **10 классов**
  - Простое обобщение с 2х!
  - Сейчас это увидим

Skillbox



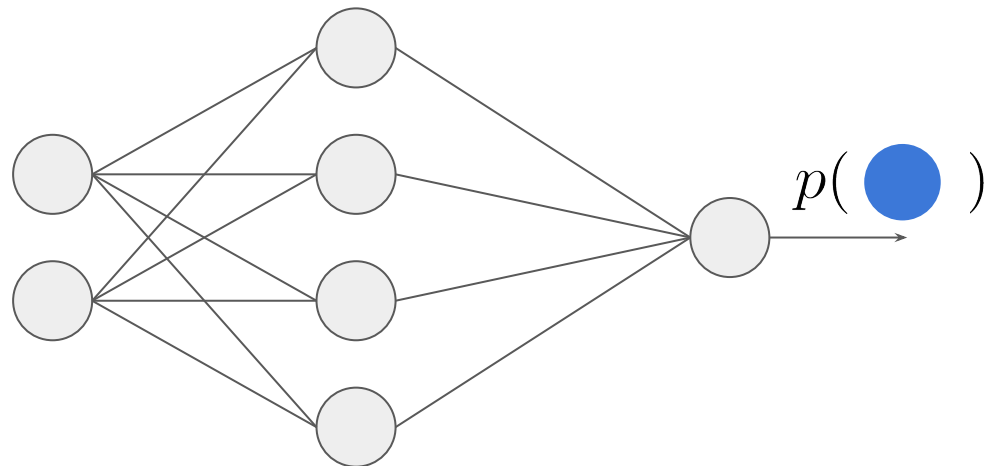
$p(\text{●})$

0.10

0.80

0.95

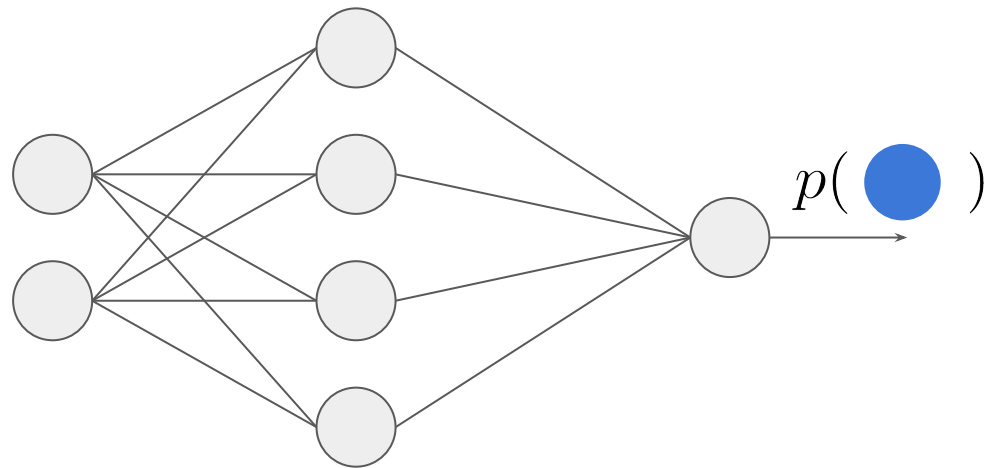
0.90



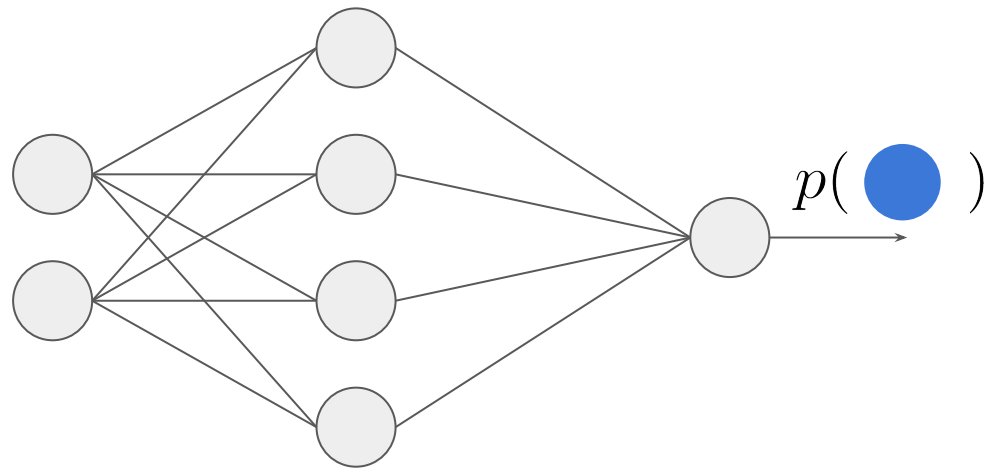


# Skillbox

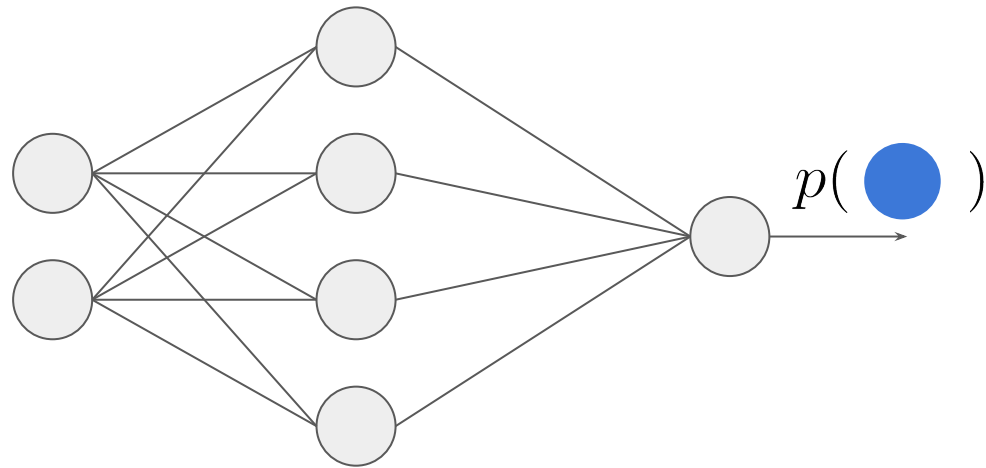
	$p(\text{red})$	$p(\text{blue})$
red	0.90	0.10
blue	0.20	0.80
blue	0.05	0.95
red	0.10	0.90



	$p(\text{red})$	$p(\text{blue})$
red	0.90	0.10
blue	0.20	0.80
blue	0.05	0.95
red	0.10	0.90

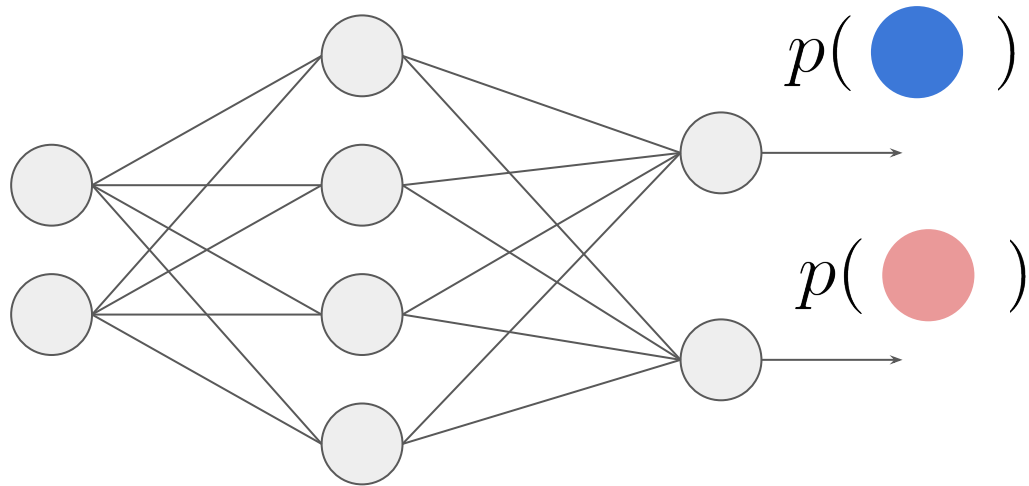


	$p(\text{red})$	$p(\text{blue})$
red	0.90	0.10
blue	0.20	0.80
blue	0.05	0.95
red	0.10	0.90

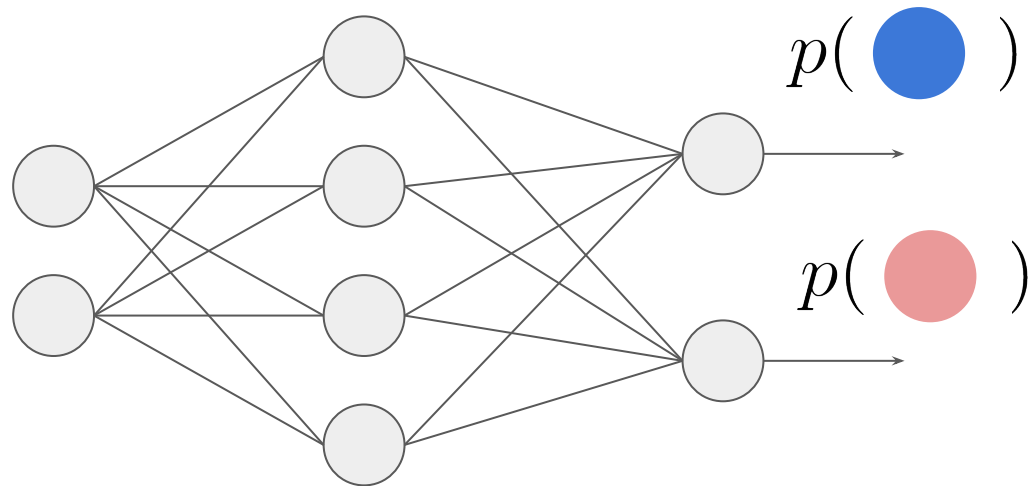


*Можно сделать два выхода?*

	$p(\text{red})$	$p(\text{blue})$
red	0.90	0.10
blue	0.20	0.80
blue	0.05	0.95
red	0.10	0.90



	$p(\text{red})$	$p(\text{blue})$
red	0.90	0.10
blue	0.20	0.80
blue	0.05	0.95
red	0.10	0.90



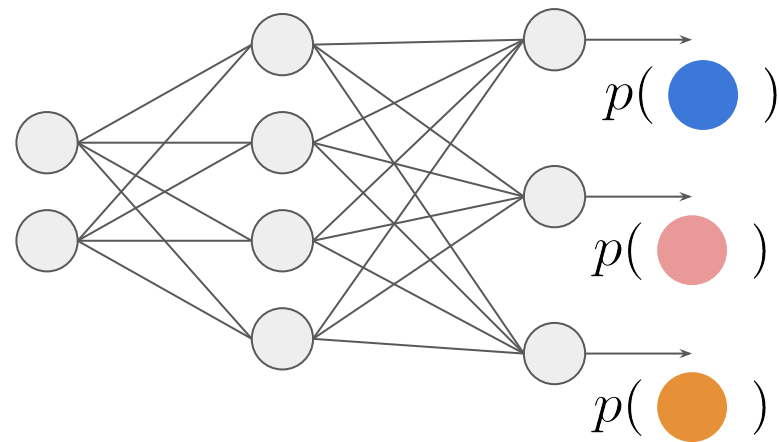
*А если три класса?*

$p(\text{red})$   $p(\text{blue})$   $p(\text{orange})$

red	0.90	0.05	0.05
blue	0.10	0.80	0.10
orange	0.05	0.92	0.03
red	0.70	0.20	0.10

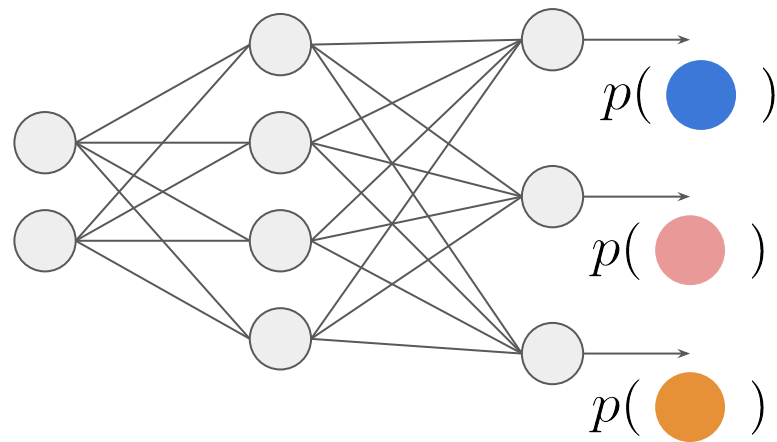
*А если три класса?*

	$p(\text{red})$	$p(\text{blue})$	$p(\text{orange})$
red	0.90	0.05	0.05
blue	0.10	0.80	0.10
orange	0.05	0.92	0.03
red	0.70	0.20	0.10



*А если три класса?*

	$p(\text{red})$	$p(\text{blue})$	$p(\text{orange})$
red	0.90	0.05	0.05
blue	0.10	0.80	0.10
orange	0.05	0.92	0.03
red	0.70	0.20	0.10



*А как считать ошибку?*



$p(\text{red})$   $p(\text{blue})$   $p(\text{orange})$

red	0.90	0.05	0.05
blue	0.10	0.80	0.10
orange	0.05	0.92	0.03
red	0.70	0.20	0.10

*А как считать ошибку?*

	$p(\text{red})$	$p(\text{blue})$	$p(\text{orange})$
red	0.90	0.05	0.05
blue	0.10	0.80	0.10
orange	0.05	0.92	0.03
red	0.70	0.20	0.10

Выбираем вероятность правильного класса

Это **НЕ** обязательно максимум в строчке -- алгоритм мог ошибиться

*А как считать ошибку?*

	$p(\text{red})$	$p(\text{blue})$	$p(\text{orange})$	$\hat{p}$
red		0.05	0.05	0.90
blue	0.10		0.10	0.80
orange	0.05	0.92		0.03
red		0.20	0.10	0.70

*А как считать ошибку?*

	$p(\text{red})$	$p(\text{blue})$	$p(\text{orange})$	$\hat{p}$
red		0.05	0.05	0.90
blue	0.10		0.10	0.80
orange	0.05	0.92		0.03
red		0.20	0.10	0.70

$$NLL(\mathbf{P}, \mathbf{y}) = -\frac{1}{N} \sum_{i=1}^N \log(\hat{p}(\mathbf{p}_i, y_i))$$

*А как считать ошибку?*

	$p(\text{red})$	$p(\text{blue})$	$p(\text{orange})$	$\hat{p}$
red	0.90	0.05	0.05	0.90
blue	0.10	0.80	0.10	0.80
orange	0.05	0.92	0.03	0.03
red	0.70	0.20	0.10	0.70

$$NLL(\mathbf{P}, \mathbf{y}) = -\frac{1}{N} \sum_{i=1}^N \log(\hat{p}(\mathbf{p}_i, y_i))$$

А как считать ошибку?

	$p(\text{red})$	$p(\text{blue})$	$p(\text{orange})$	$\hat{p}$
red	0.90	0.05	0.05	0.90
blue	0.10	0.80	0.10	0.80
orange	0.05	0.92	0.03	0.03
red	0.70	0.20	0.10	0.70

$$NLL(\mathbf{P}, \mathbf{y}) = -\frac{1}{N} \sum_{i=1}^N \log(\hat{p}(\mathbf{p}_i, y_i))$$

А как считать ошибку?

	$p(\text{red})$	$p(\text{blue})$	$p(\text{orange})$	$\hat{p}$
red	0.90	0.05	0.05	0.90
blue	0.10	0.80	0.10	0.80
orange	0.05	0.92	0.03	0.03
red	0.70	0.20	0.10	0.70

$$NLL(\mathbf{P}, \mathbf{y}) = -\frac{1}{N} \sum_{i=1}^N \log(\hat{p}(\mathbf{p}_i, y_i))$$

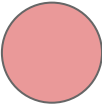


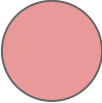
*А как считать ошибку? Аналогично случаю двух классов*

	$p(\text{red})$	$p(\text{blue})$	$p(\text{orange})$	$\hat{p}$
red	0.90	0.05	0.05	0.90
blue	0.10	0.80	0.10	0.80
orange	0.05	0.92	0.03	0.03
red	0.70	0.20	0.10	0.70

$$NLL(\mathbf{P}, \mathbf{y}) = -\frac{1}{N} \sum_{i=1}^N \log(\hat{p}(\mathbf{p}_i, y_i))$$

*Какие условия на выходы сети?*



	$p(\text{red})$	$p(\text{blue})$	$p(\text{orange})$	$\hat{p}$
	0.90	0.05	0.05	0.90
	0.10	0.80	0.10	0.80
	0.05	0.92	0.03	0.03
	0.70	0.20	0.10	0.70

$$NLL(\mathbf{P}, \mathbf{y}) = -\frac{1}{N} \sum_{i=1}^N \log(\hat{p}(\mathbf{p}_i, y_i))$$

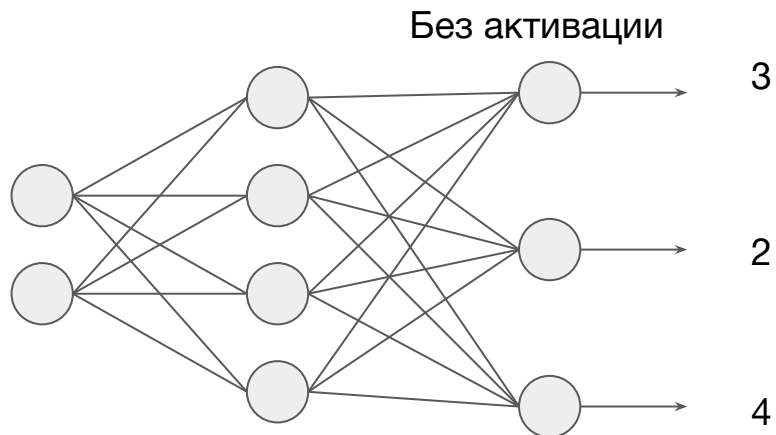
*Какие условия на выходы сети? В сумме должна быть единица. Нужна нормировка ...*

# SoftMax

$$p_i^k = \frac{\exp(\text{logit}_i^k)}{\sum_{j=1}^m \exp(\text{logit}_i^j)}$$

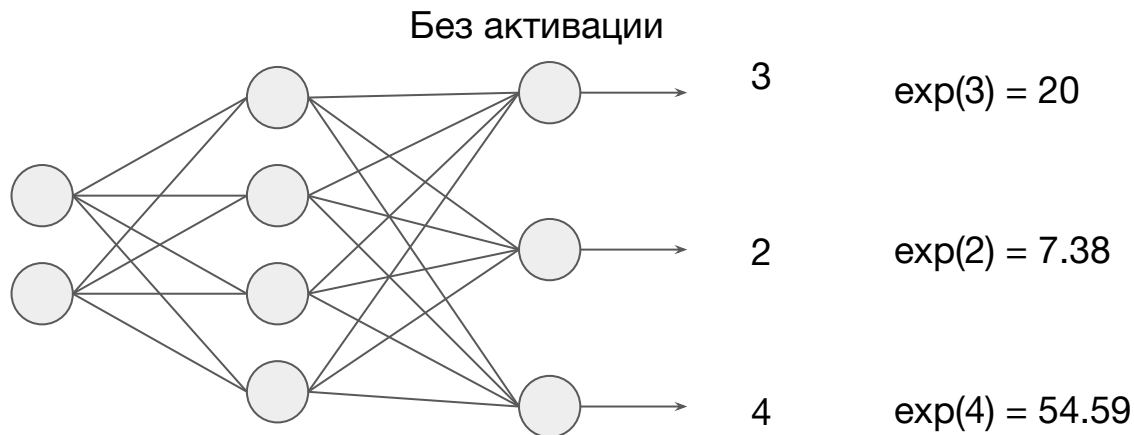
# SoftMax

$$p_i^k = \frac{\exp(\text{logit}_i^k)}{\sum_{j=1}^m \exp(\text{logit}_i^j)}$$



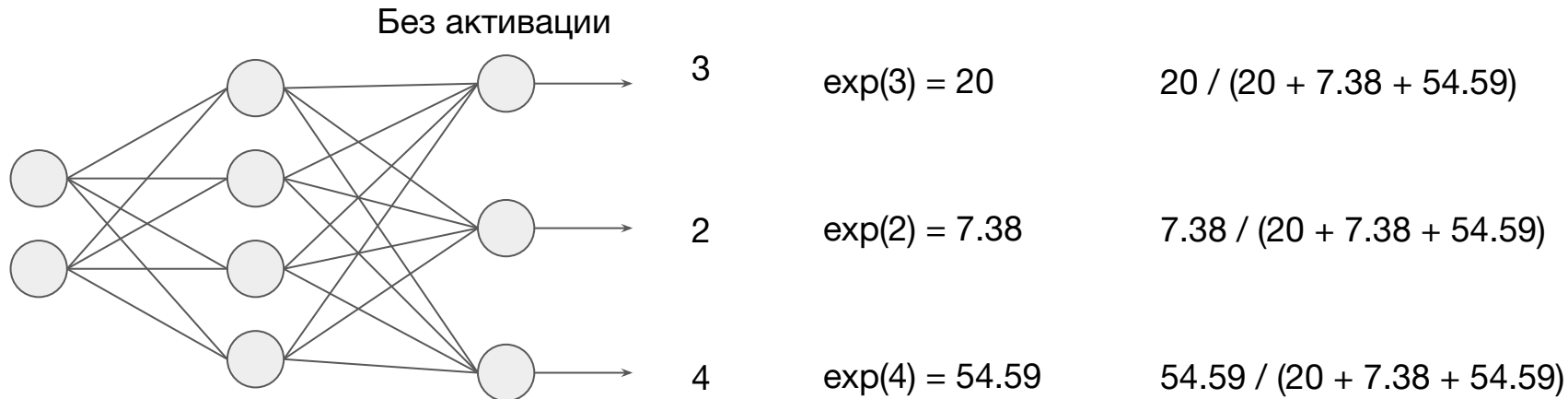
# SoftMax

$$p_i^k = \frac{\exp(\text{logit}_i^k)}{\sum_{j=1}^m \exp(\text{logit}_i^j)}$$



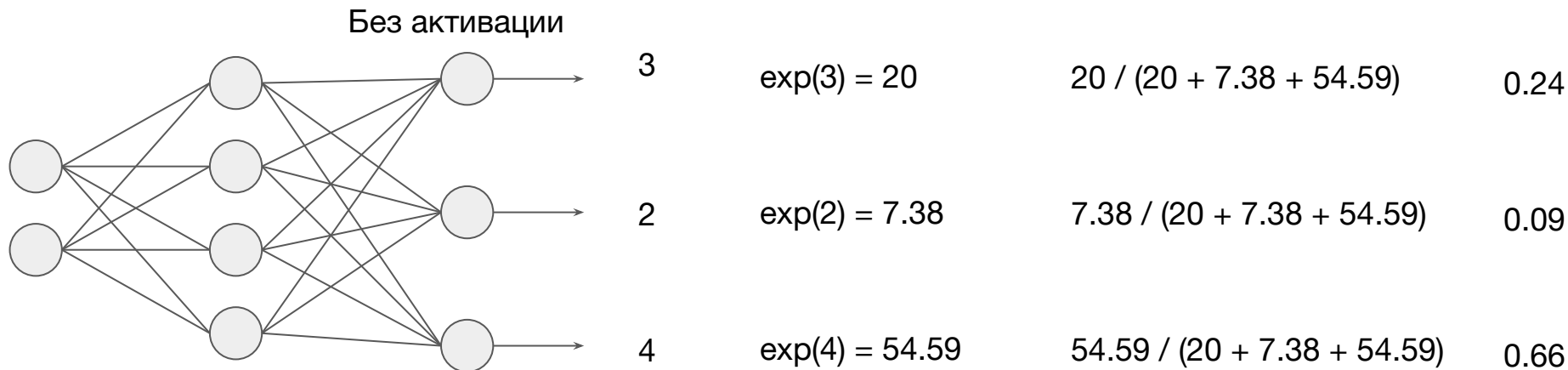
# SoftMax

$$p_i^k = \frac{\exp(\text{logit}_i^k)}{\sum_{j=1}^m \exp(\text{logit}_i^j)}$$



# SoftMax

$$p_i^k = \frac{\exp(\text{logit}_i^k)}{\sum_{j=1}^m \exp(\text{logit}_i^j)}$$



# Итог

- Для мульти классовой классификации вместо сигмоиды используется SoftMax
- Он выполняет функцию нормировки выходов сети
  - Положительные
  - В сумме дают 1
- В качестве функции потерь используется известный нам логлосс:
  - Categorical cross entropy