

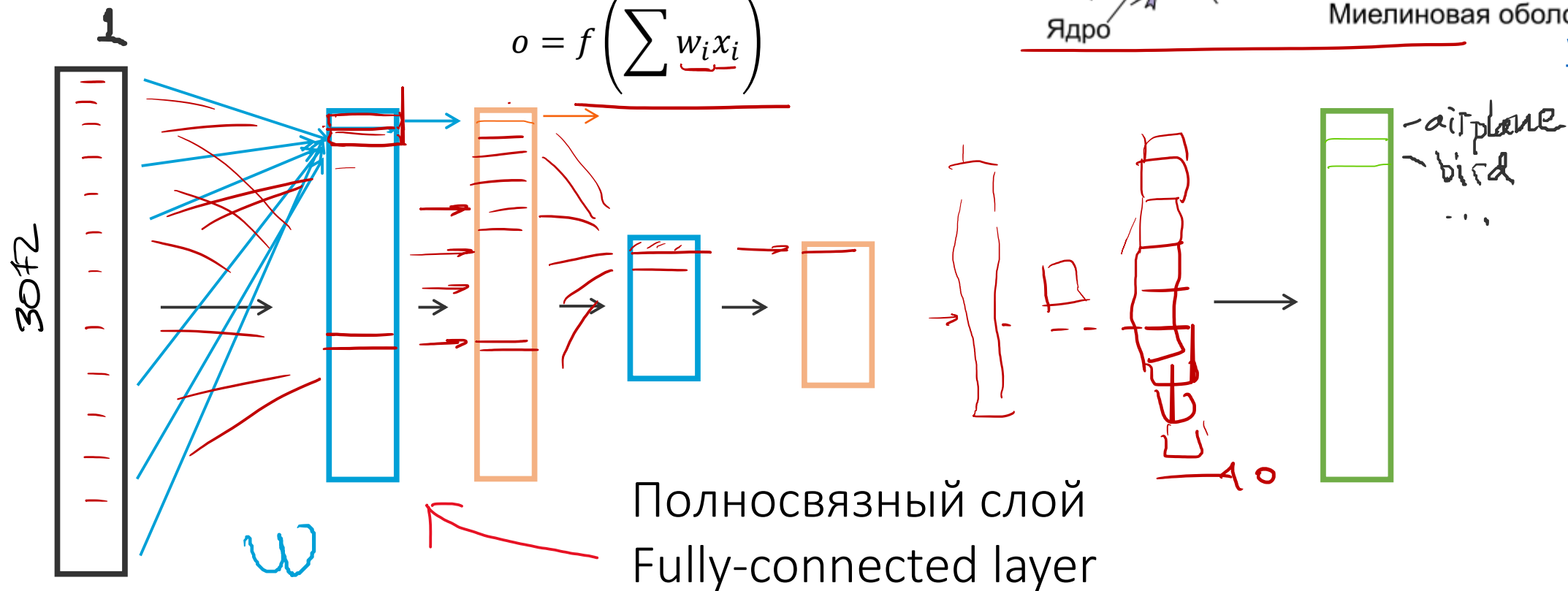
Лекция 4

Нейронные сети —
подробности

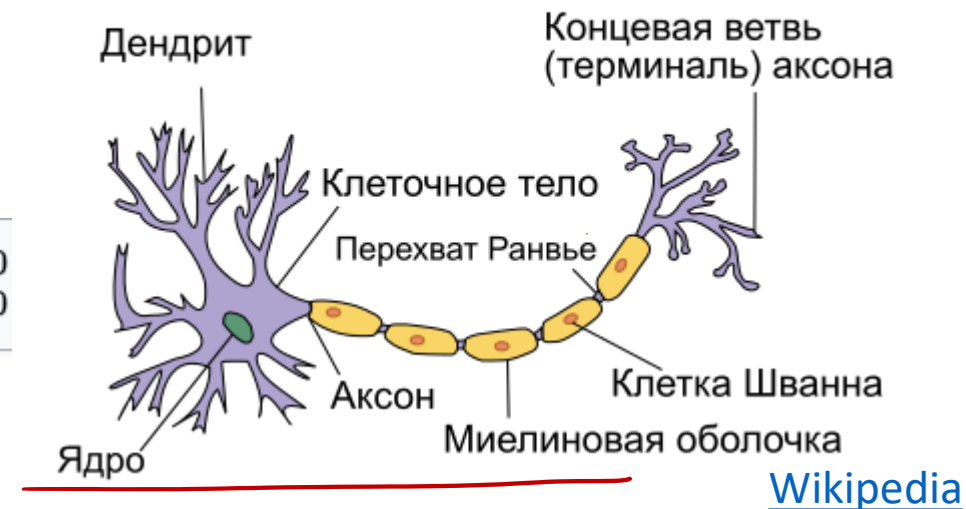
Нейронная сеть Neural network



$$o = f\left(\sum w_i x_i\right)$$

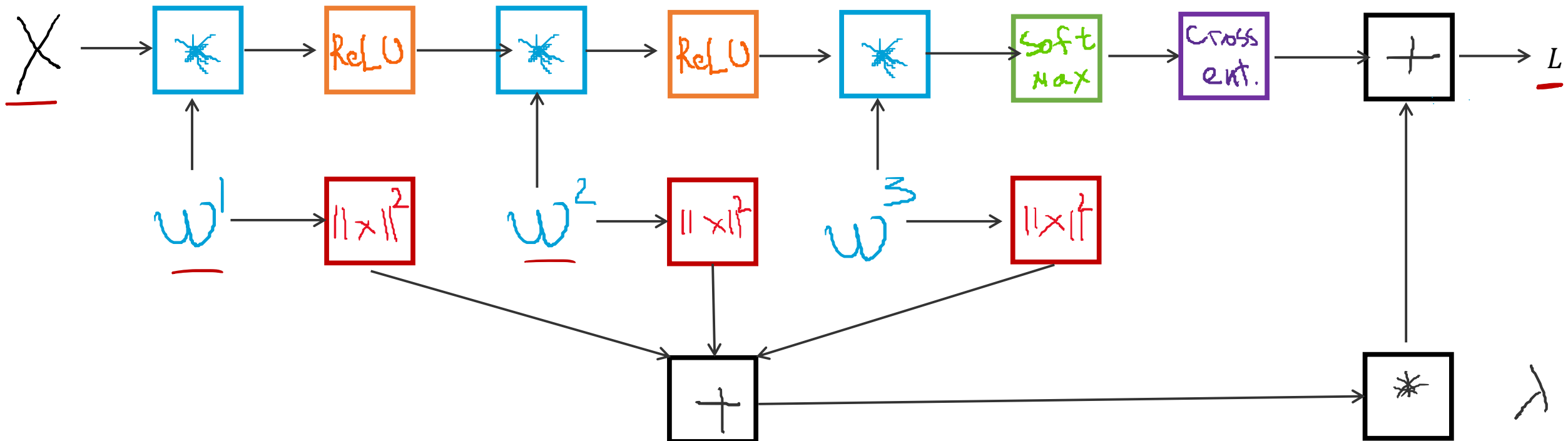


Типичная структура нейрона



Граф вычислений Computational graph

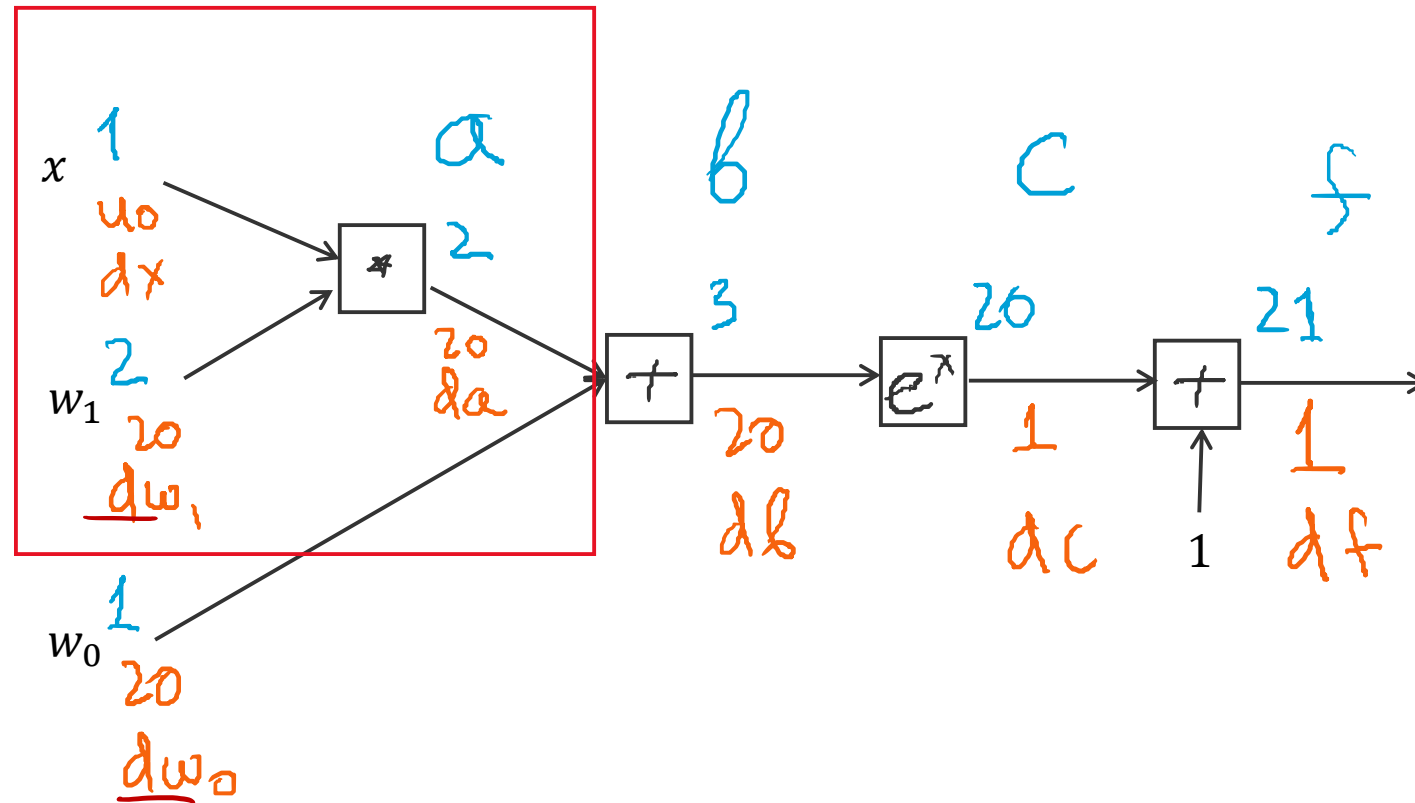
$$L = - \sum_j \ln p(c = y_j | x_j) + \lambda R(\omega)$$



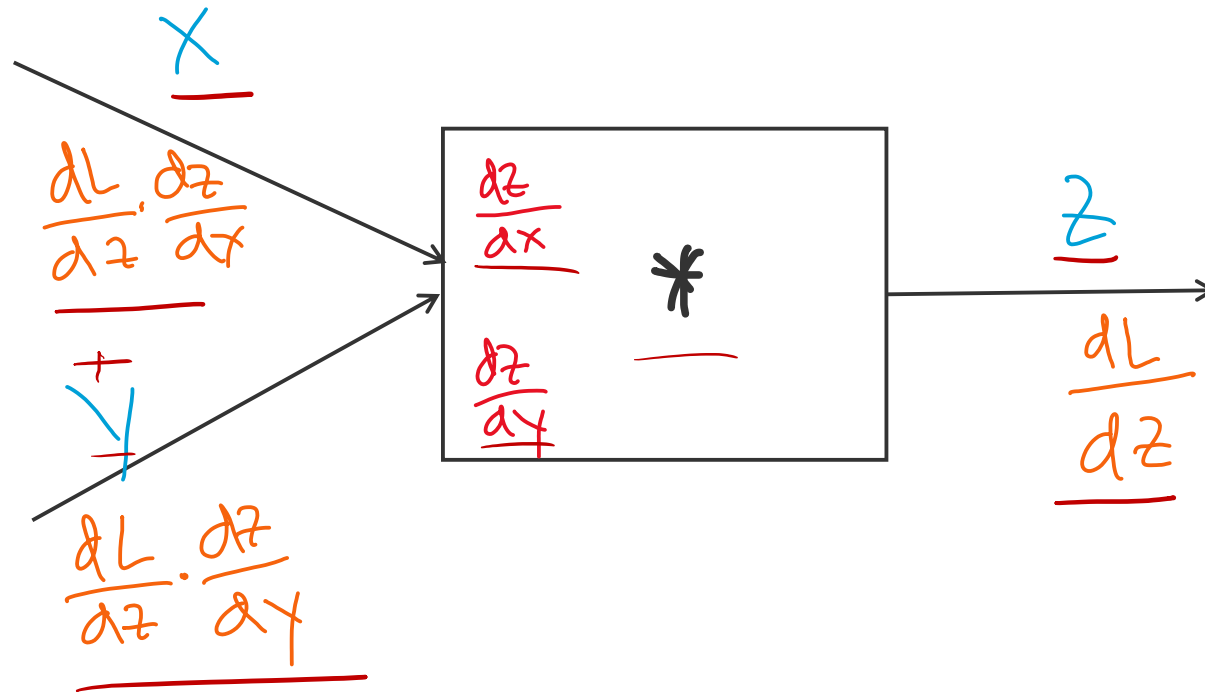
Обратное распространение ошибки Backpropagation

$$f(g(x)) \quad \frac{df}{dx} = \frac{df}{dg} \frac{dg}{dx}$$

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



Общая схема вычисления градиента



...

А ТЕПЕРЬ ТО ЖЕ САМОЕ



С МАТРИЦАМИ

А теперь с матрицами

$$f = \frac{1}{2} \|X \cdot W\|_2^2$$

L2

X $\begin{bmatrix} 2 \\ 6 \end{bmatrix}$ $\frac{2 \times 1}{RC}$ really cool a $\begin{bmatrix} 24 \\ 00 \end{bmatrix}$

$$\nabla_x f$$

$\begin{bmatrix} 1 & 2 \end{bmatrix}$

W

$$b = \sum a_{ij}^2$$

$$\frac{df}{da_{ij}} = \frac{db}{db} \cdot \frac{db}{da_{ij}} = a_{ij}$$

0.5 2a_{ij}

$$\nabla_a f = \nabla_a$$

$$\|X\|_2^2$$

b
20

0.5
db

$$\nabla_a f = a$$

1/2

$$\nabla_a f = a$$

f
10

1
df

А теперь с матрицами

$$f = \frac{1}{2} \|X \cdot W\|_2^2$$

Handwritten notes and matrices:

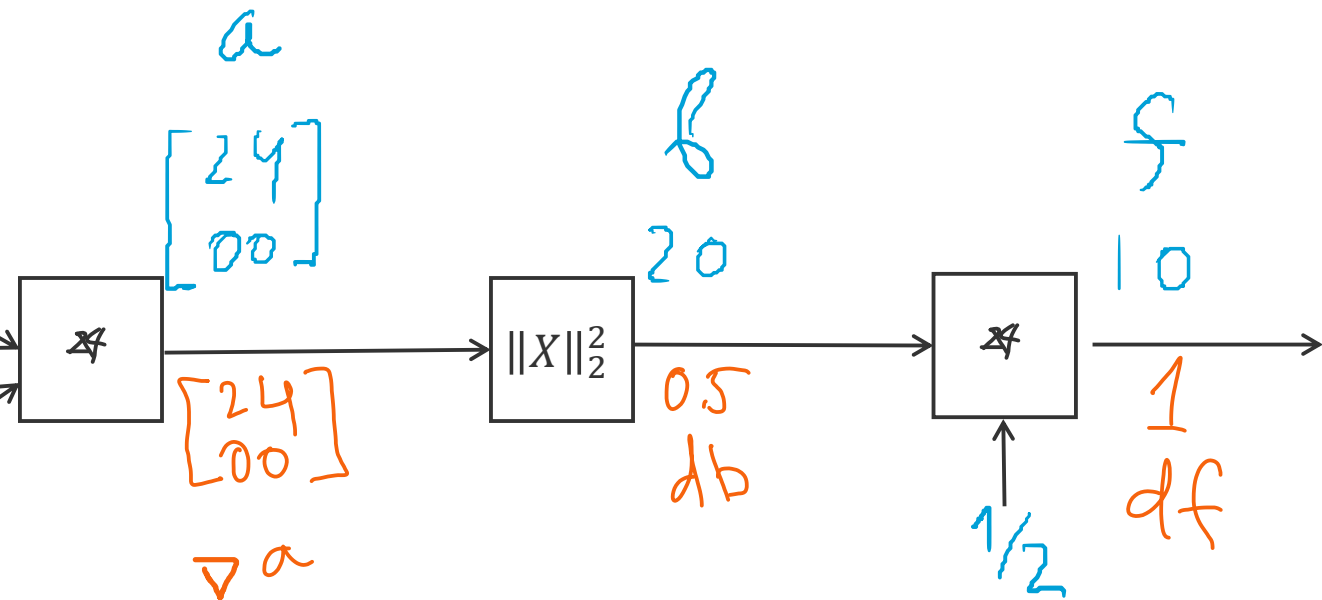
$\nabla_a f = X^T \cdot \nabla_a f$
Dimensions: 1×2 (row vector) \cdot 2×2 (matrix) \rightarrow 1×2 (row vector)

$\nabla_a f$
Dimensions: 2×2 (matrix)

$a = X \cdot w$
Dimensions: 2×2 (matrix) \cdot 2×1 (column vector)

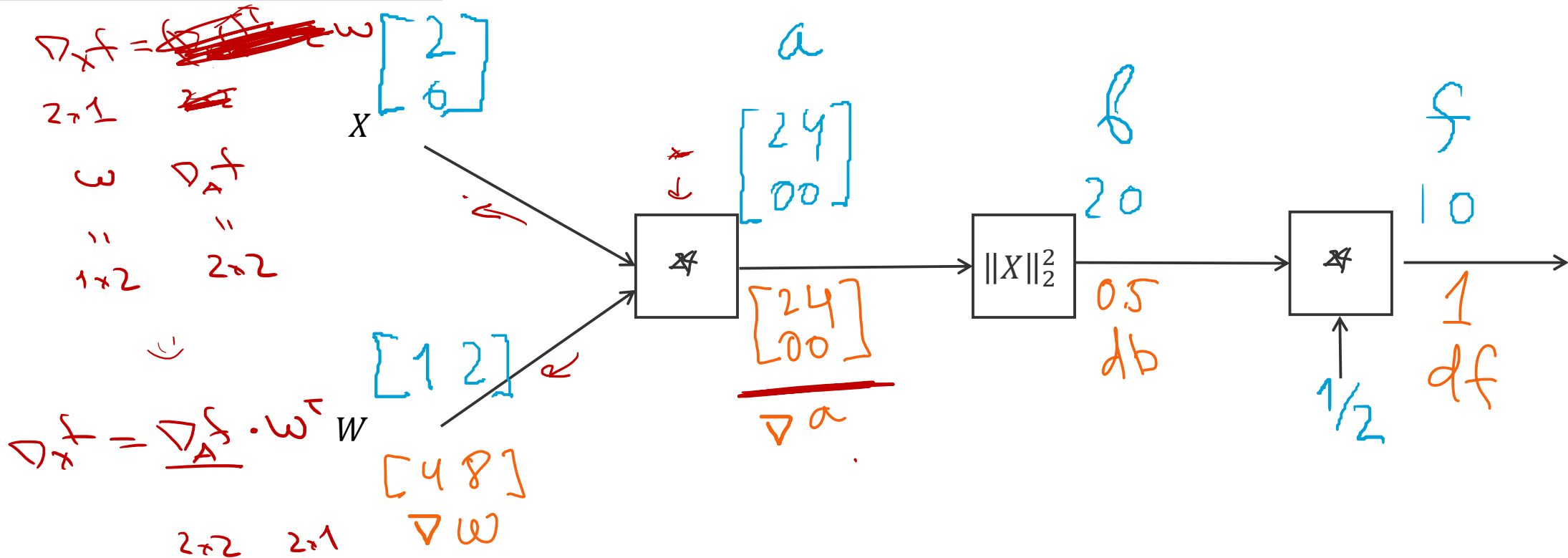
Matrix X : $\begin{bmatrix} 2 \\ 6 \end{bmatrix}$

Matrix W : $\begin{bmatrix} 1 & 2 \end{bmatrix}$
Dimensions: 1×2



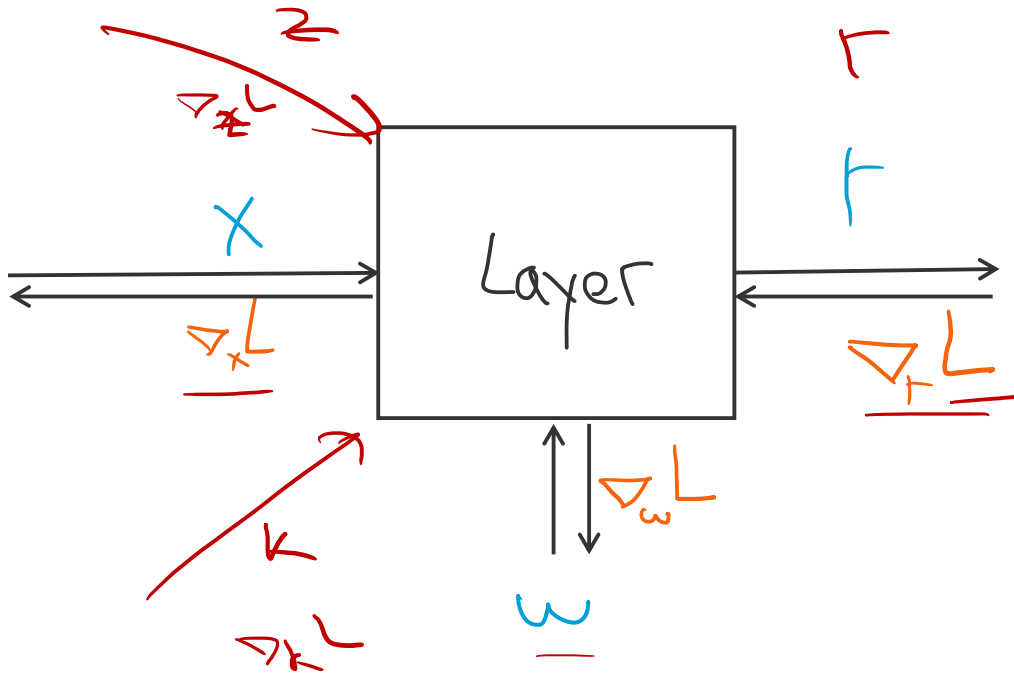
А теперь с матрицами

$$f = \frac{1}{2} \|\underline{X} \cdot \underline{W}\|_2^2$$



Интерфейс слоя

Layer interface

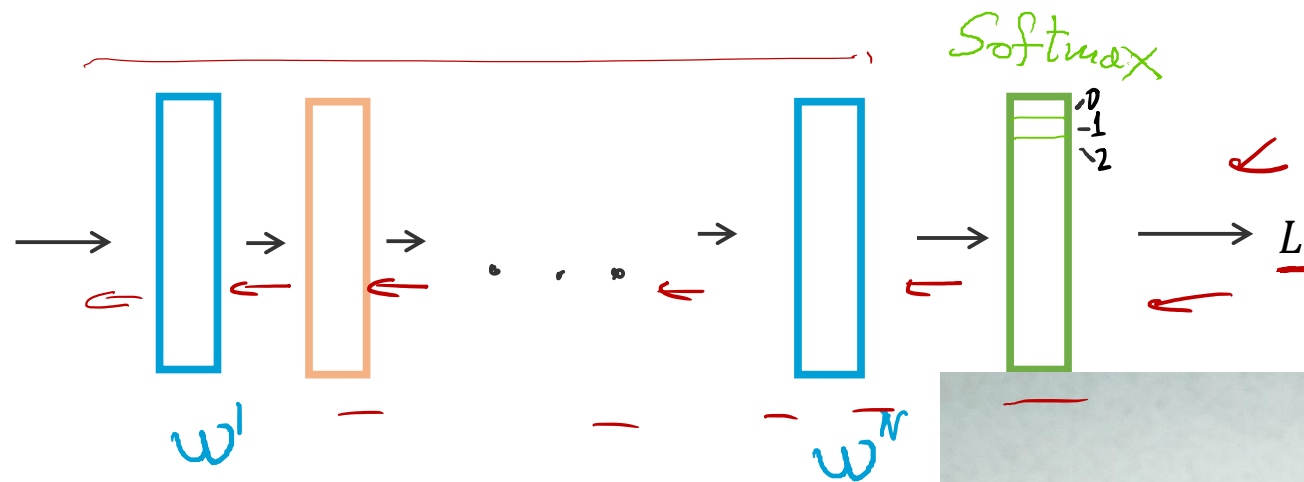
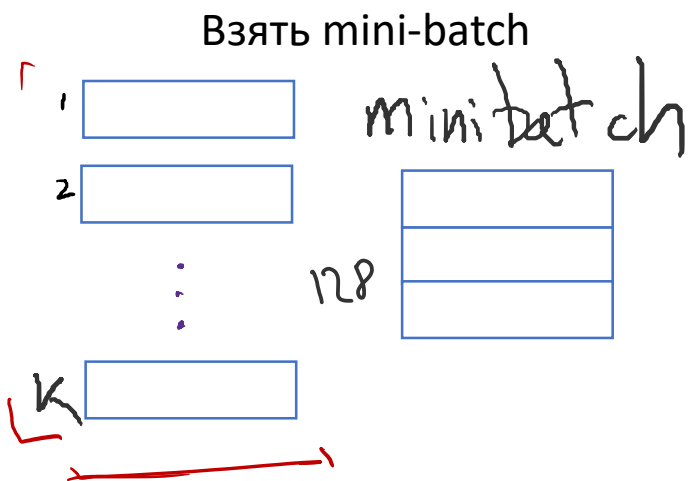


```
class Layer:
    def forward(self, x):
        result = ... # промежуточные вычисления
        self.x = x # сохраняем значения, которые нам
                    # понадобятся при обратном проходе
        return result

    def backward(self, grad):
        dx = ... # используем сохраненные значения чтобы
        dw = ... # вычислить градиент по x и по w
        self.w.grad += dw # аккумулируем градиент dw
        return dx
```

Общая схема тренировки

Прямой проход (forward pass) – посчитать loss



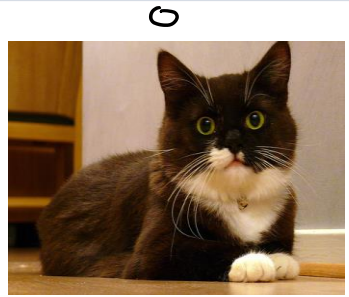
Обратный проход (backward pass) – пос

Обновить параметры

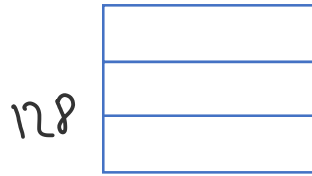
$$\begin{aligned}\overline{w}^1 &= \overline{w}^1 - \eta \overline{\nabla}_{w^1} L \\ &\dots \\ \overline{w}^n &= \overline{w}^n - \eta \overline{\nabla}_{w^n} L\end{aligned}$$



Multi-class classification



mini batch



w^1



...



w^N



$\frac{1}{2}$

→

L

Softmax

$$p(C = 0|x) = \frac{e^{y_0}}{e^{y_0} + e^{y_1} + \dots + e^{y_n}} = \frac{e^{y_0}}{\sum_i e^{y_i}}$$

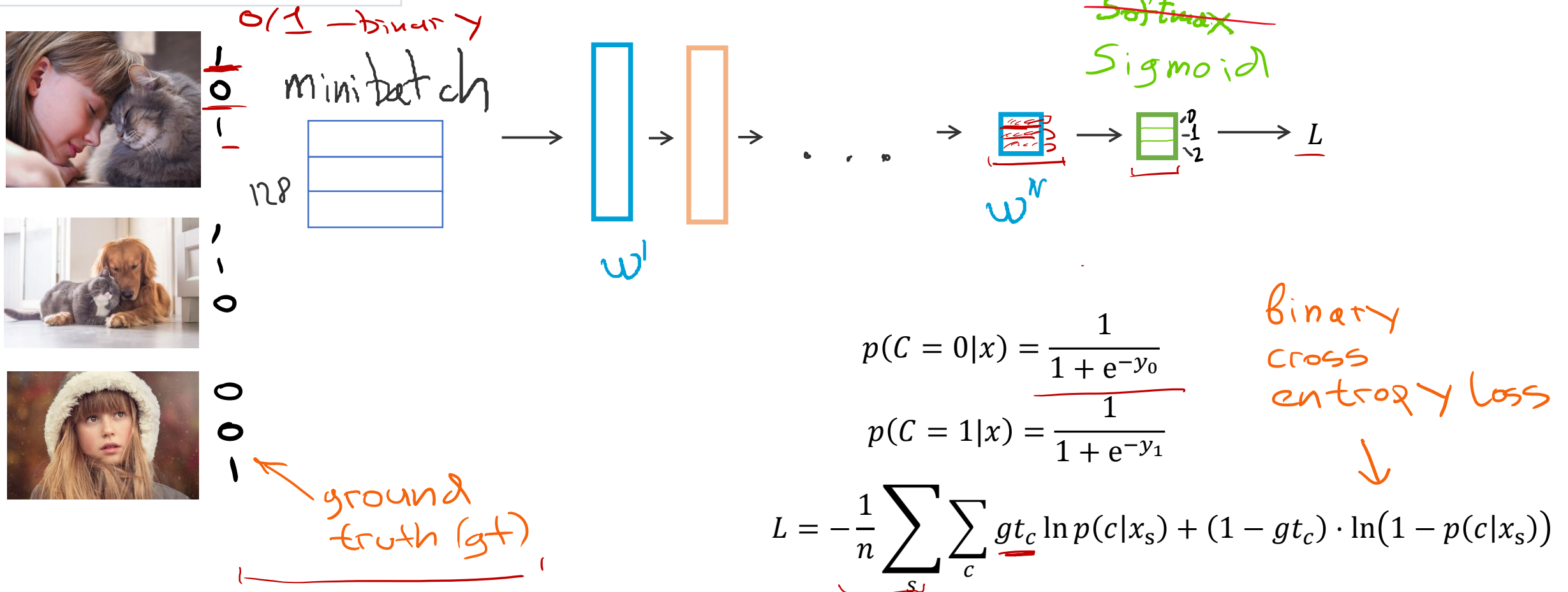
...

cross-entropy

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

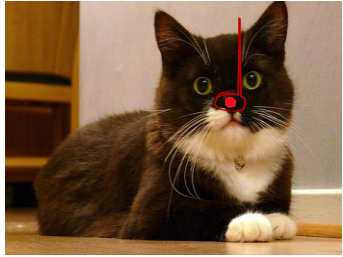
ground truth (gt)

Multi-class labeling



Regression

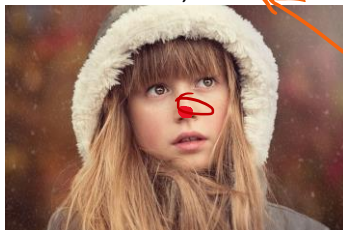
0.7, 0.3



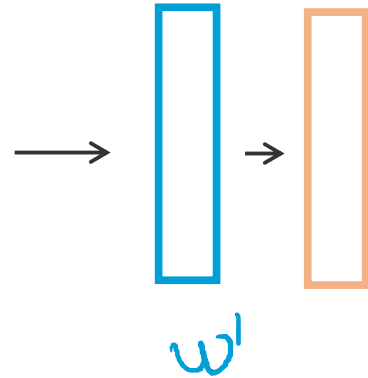
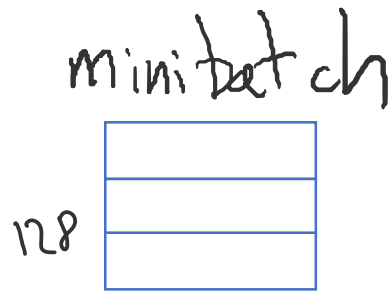
0.6, 0.5



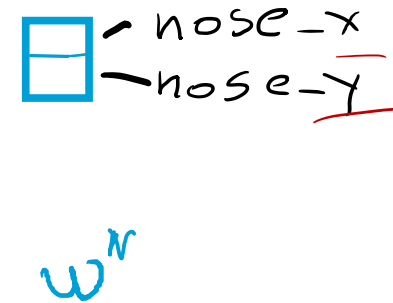
0.5, 0.5



ground
truth (gt)



...



→ L

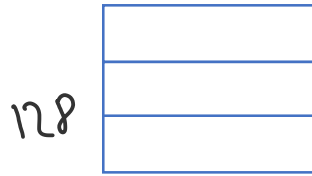
L1 loss
L2 loss

$$L = \frac{1}{n} \sum_s \sum_i |y_i - \underline{gt_i}|^2$$

Regression



mini batch



128



w^1



...



w^r

nose-x
nose-y
eye-x
eye-y
...



L

L1 loss

$$L = \frac{1}{n} \sum_s \sum_i |y_i - gt_i|$$

Восход солнца вручную



Владимир Куш "Сизиф"
Экфрасис

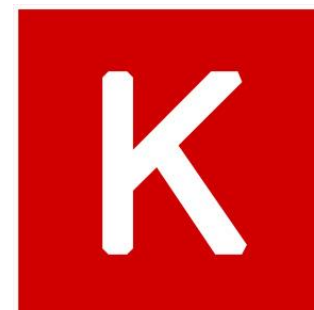
Библиотеки для глубокого обучения! Deep Learning Frameworks!



Caffe

theano

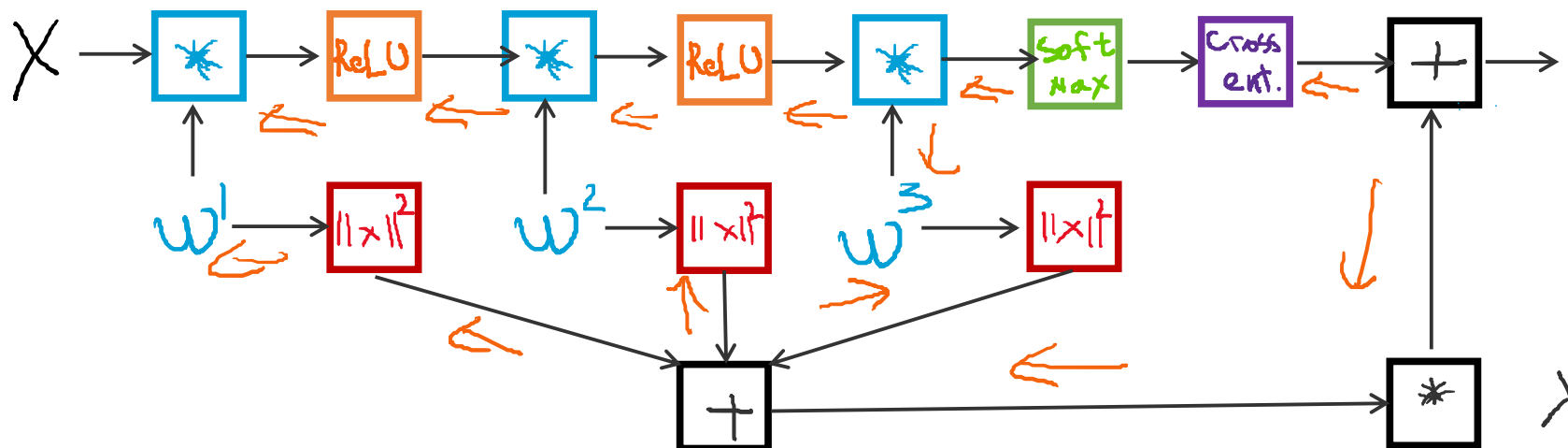
dmlc
mxnet



PYTORCH

Что они позволяют

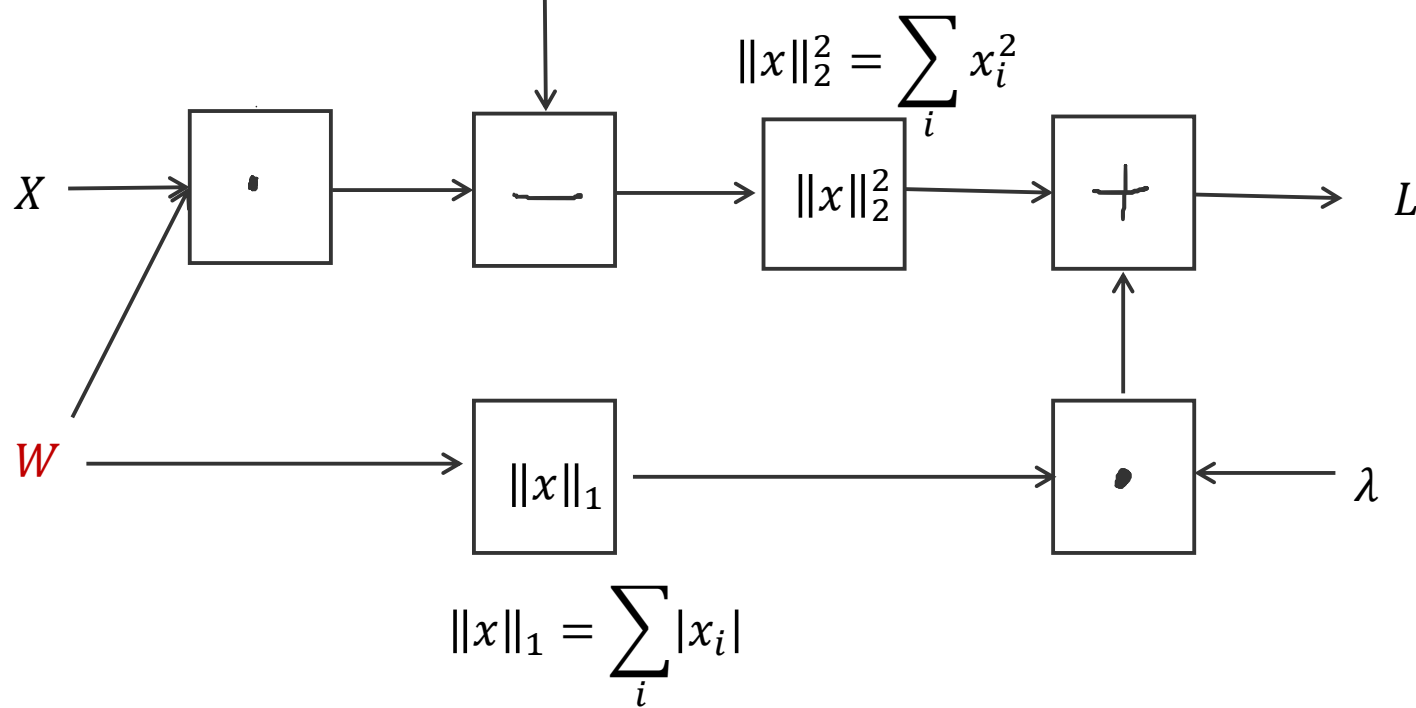
1. Задать граф вычислений
2. Посчитать градиенты на графе через backpropagation
3. Сделать это все на GPU



PyTorch: Tensors

2

$$L = \| \underbrace{X \cdot W - y}_{y} \|_2^2 + \lambda \| W \|_1$$



```
np.random.seed(2)
Y = np.random.randn(2, 5)
X = np.random.randn(2, 3)
W = np.random.randn(3, 5)
lam = 0.5

pred = X.dot(W)
pred_loss = np.sum((pred - Y)**2)
reg_loss = lam * np.sum(np.abs(W))
loss = pred_loss + reg_loss
```

```
import torch
Y = torch.Tensor(np.random.randn(2, 5))
X = torch.Tensor(np.random.randn(2, 3))
W = torch.Tensor(np.random.randn(3, 5))
lam = 0.5

pred = X.matmul(W)
pred_loss = torch.sum((pred - Y)**2)
reg_loss = lam * torch.sum(torch.abs(W))
loss = pred_loss + reg_loss
```

PyTorch: autograd

```
import torch
Y = torch.Tensor(np.random.randn(2, 5))
X = torch.Tensor(np.random.randn(2, 3))
W = torch.Tensor(np.random.randn(3, 5))
lam = 0.5
pred = X.matmul(W)
pred_loss = torch.sum((pred - Y)**2)
reg_loss = lam * torch.sum(torch.abs(W))
loss = pred_loss + reg_loss
```

forward

```
Y = torch.Tensor(np.random.randn(2, 5)).cuda()
X = torch.Tensor(np.random.randn(2, 3)).cuda()
W = torch.Tensor(np.random.randn(3, 5)).cuda().requires_grad_()
l = 0.5

for i in range(100):
    pred = X.matmul(W)
    pred_loss = torch.sum((pred - Y)**2)
    reg_loss = l * torch.sum(torch.abs(W))
    loss = pred_loss + reg_loss

    loss.backward()
    W.data.add_(-0.1*W.grad.data)
    W.grad.zero_()
```

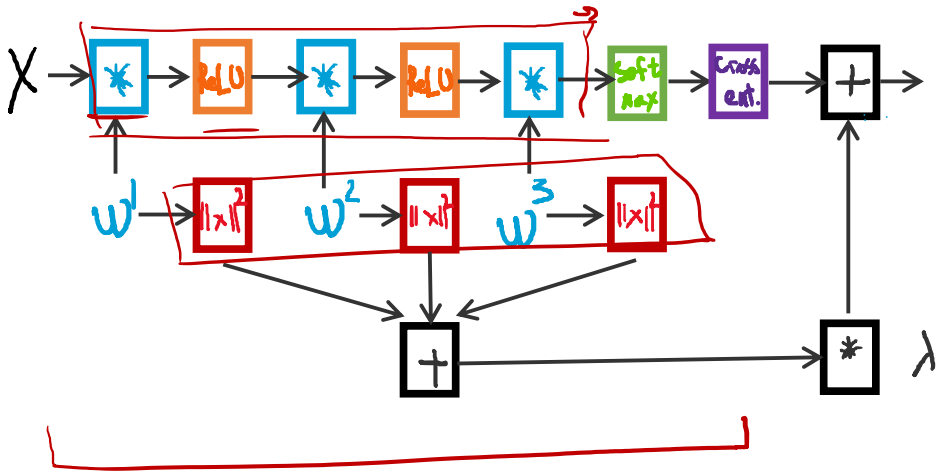
BUT WAIT



THERE'S MORE

memegenerator.net

PyTorch: Modules



```
import torch.optim as optim
y = torch.LongTensor(np.random.randint(0,3,size=2))
```

```
nn = torch.nn.Sequential(
    torch.nn.Linear(3, 10),
    torch.nn.ReLU(),
    torch.nn.Linear(10, 20),
    torch.nn.ReLU(),
    torch.nn.Linear(20, 3)
```

```
)
optimizer = optim.SGD(nn.parameters(), lr=0.01, weight_decay=0.05)
```

```
for i in range(100):
    optimizer.zero_grad()
    pred = nn(X)
    criterion = torch.nn.CrossEntropyLoss()
    loss = criterion(pred, y)
    loss.backward()
    optimizer.step()
```


PyTorch: Modules

```
nn = torch.nn.Sequential(  
    torch.nn.Linear(3, 10),  
    torch.nn.ReLU(),  
    torch.nn.Linear(10, 20),  
    torch.nn.ReLU(),  
    torch.nn.Linear(20, 3))
```

```
nn(X)
```

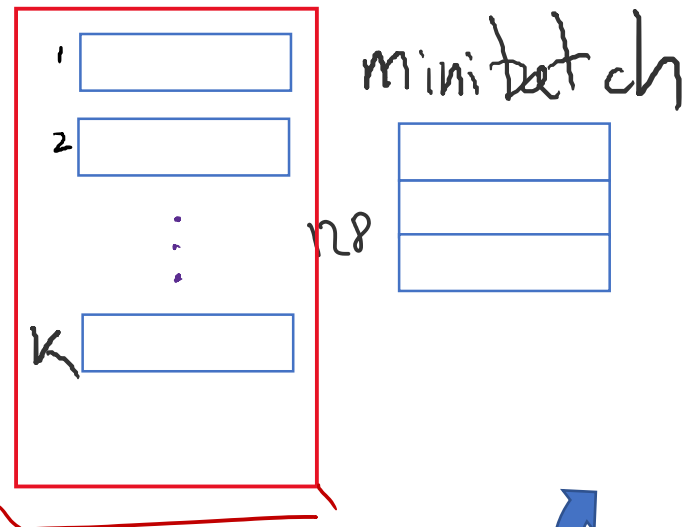
```
import torch.nn.functional as F  
  
class Net(torch.nn.Module):  
    def __init__(self):  
        super(Net, self).__init__()  
        self.fc1 = torch.nn.Linear(3, 10)  
        self.fc2 = torch.nn.Linear(10, 20)  
        self.fc3 = torch.nn.Linear(20, 3)  
  
    def forward(self, x):  
        x = F.relu(self.fc1(x))  
        x = F.relu(self.fc2(x))  
        x = self.fc3(x)  
        return x
```

```
nn = Net()  
nn(X)
```

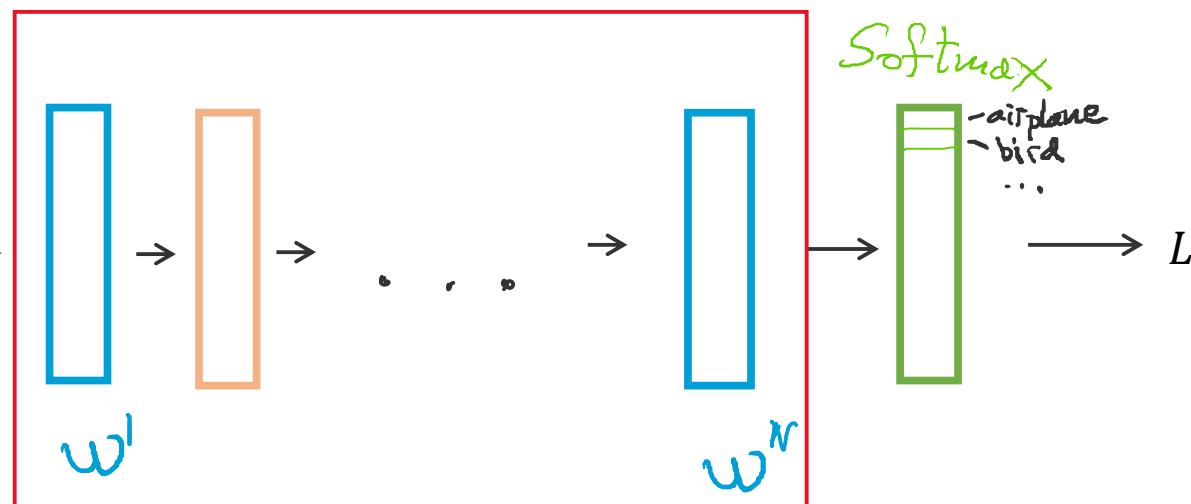


Погружаемся в детали

Взять mini-batch



Прямой проход (forward pass) – посчитать loss



Обратный проход (backward pass) – посчитать градиент

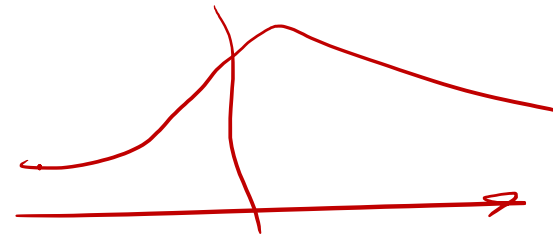
Обновить параметры

$$\overrightarrow{w^1} = \overrightarrow{w^1} - \eta \overrightarrow{\nabla_{w^1} L}$$

...

$$\overrightarrow{w^n} = \overrightarrow{w^n} - \eta \overrightarrow{\nabla_{w^n} L}$$

Подготовка данных Data preprocessing

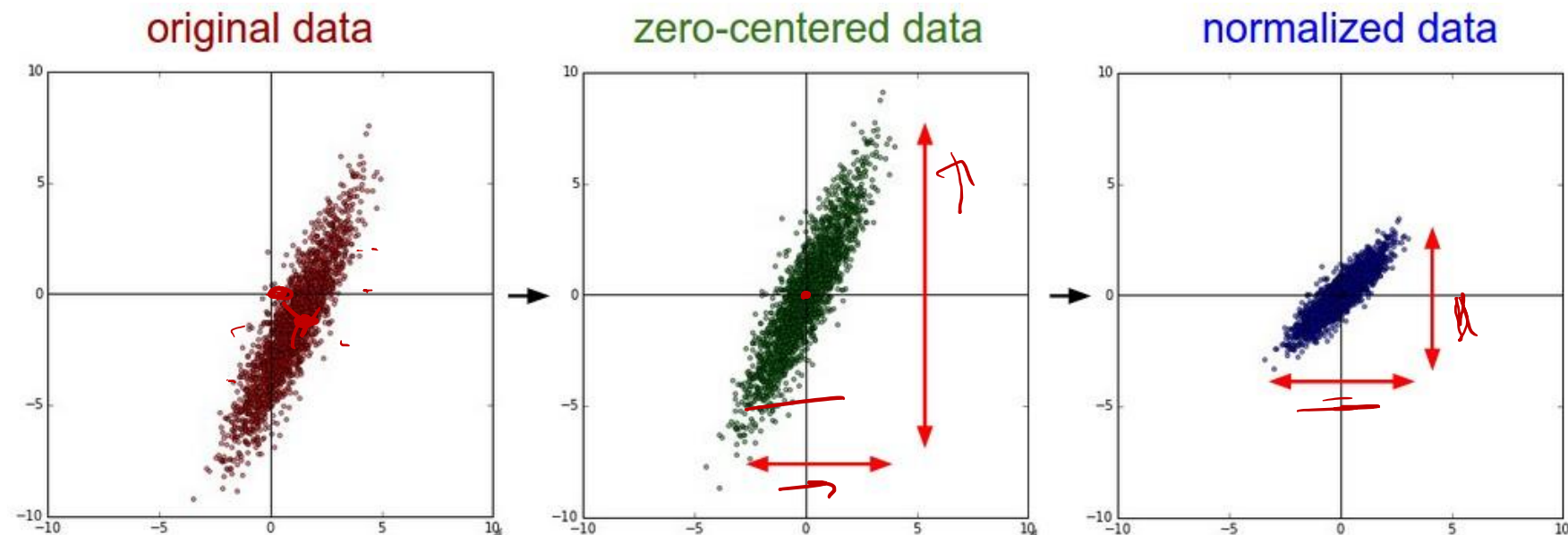


1

2

⋮

K



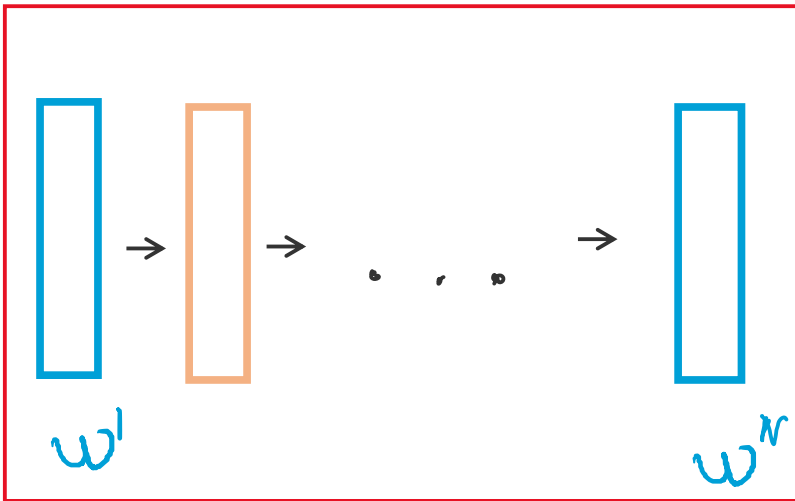
$$x' = x - \bar{x}$$

среднее

$$x' = \frac{x - \bar{x}}{\sigma}$$

$$\sigma = \sqrt{\frac{(x - \bar{x})^2}{N - 1}}$$

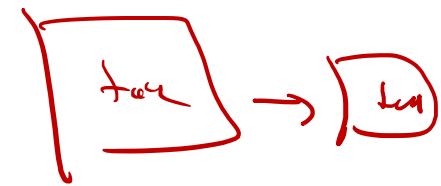
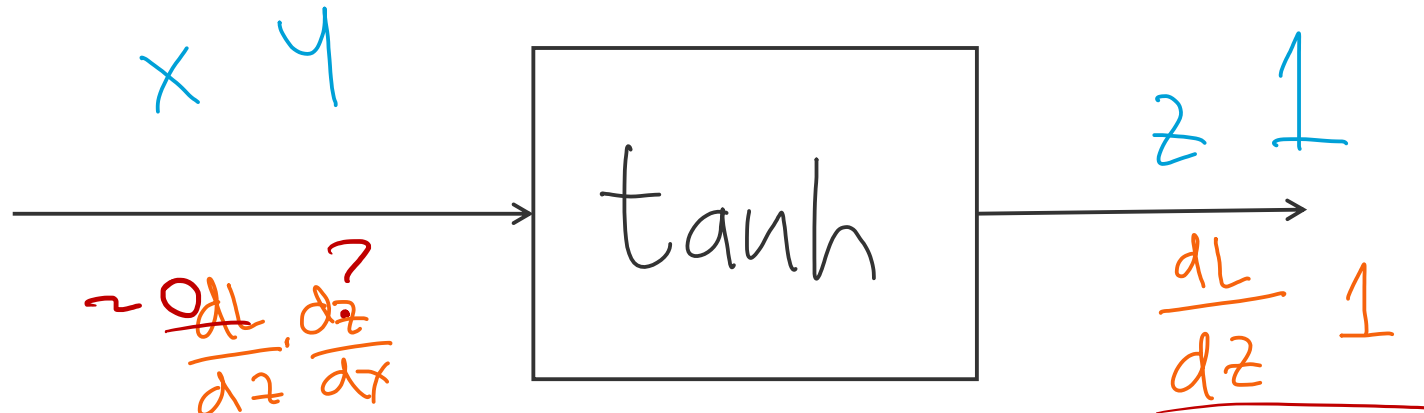
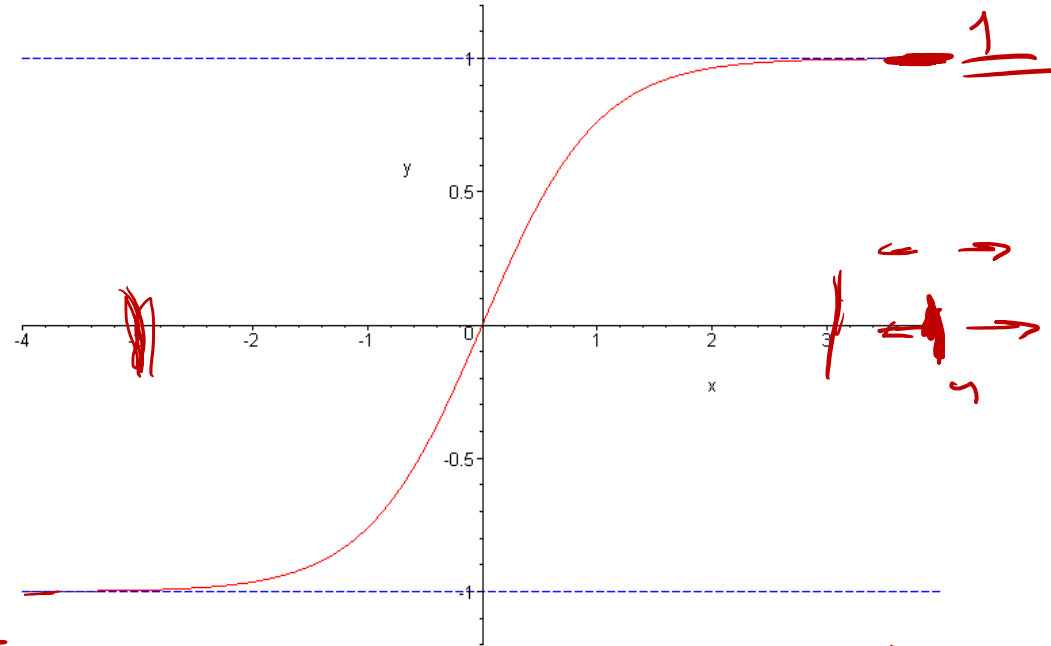
Детали работы сети



- Выбор активационной функции
- Инициализация весов

Проблема: Vanishing gradients

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

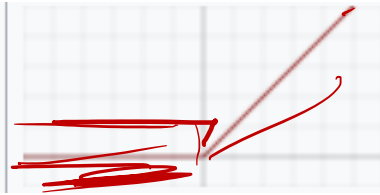


Rectifier Linear Unit

ReLU

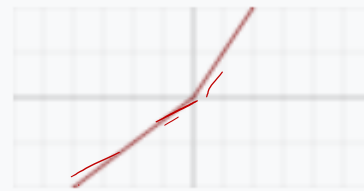
[Wikipedia](#)

Rectified linear
unit (ReLU)^[10]



$$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$$

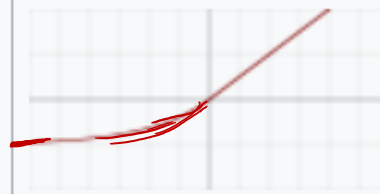
Leaky rectified
linear unit
(Leaky
ReLU)^[11]



$$f(x) = \begin{cases} 0.01x & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$$

[Rectifier Nonlinearities Improve Neural Network Acoustic Models'13](#)

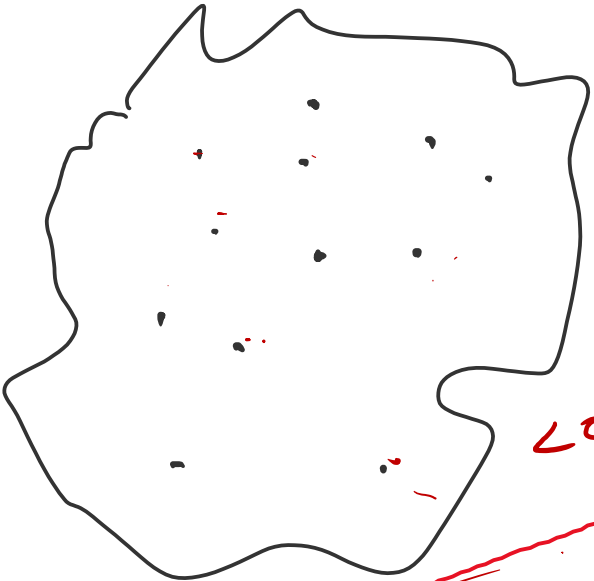
Exponential
linear unit
(ELU)^[14]



$$f(\alpha, x) = \begin{cases} \alpha(e^x - 1) & \text{for } x < 0 \\ x & \text{for } x \geq 0 \end{cases}$$

[Fast and Accurate Deep Network Learning by Exponential Linear Units \(ELUs\)'15](#)

Data 10%

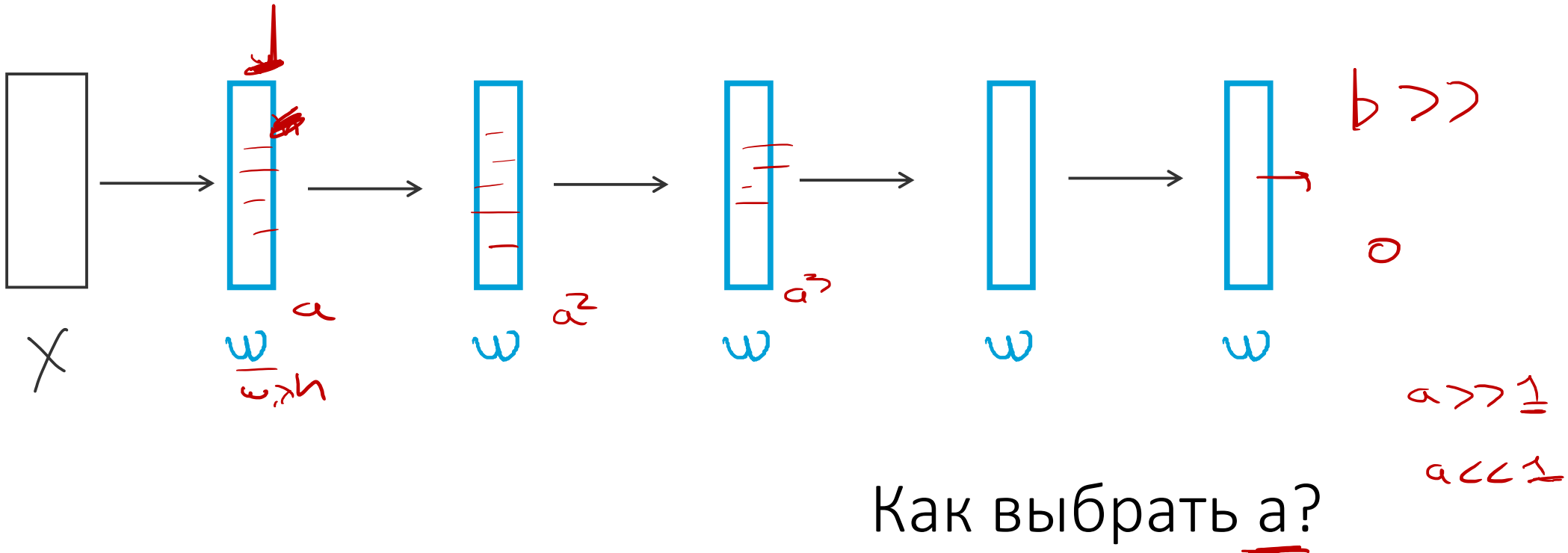


ReLU

Инициализация весов

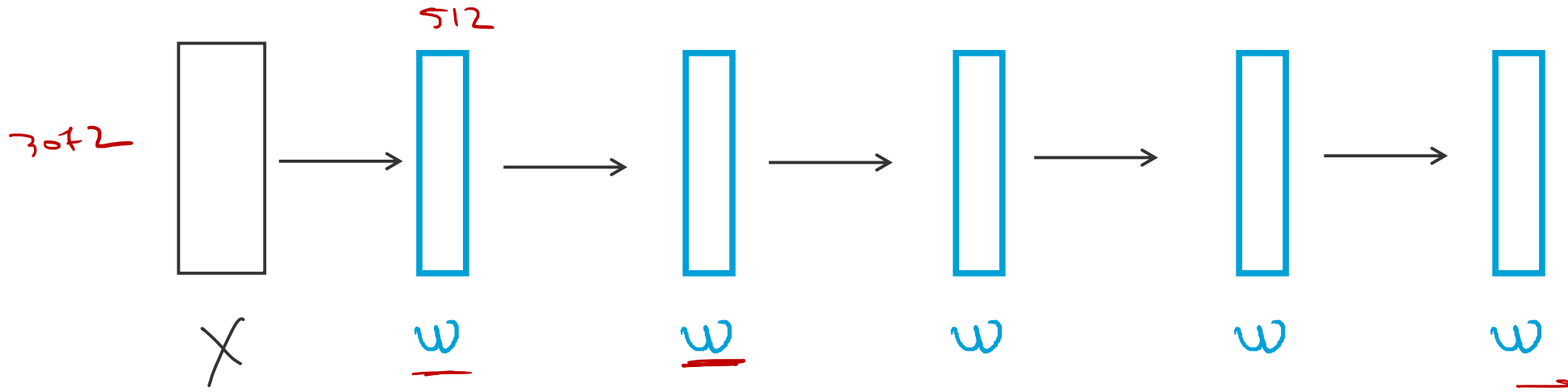
Weight initialization

$$W = \underline{a} * \underline{\text{random}(\text{width}, \text{height})}$$



Xavier initialization

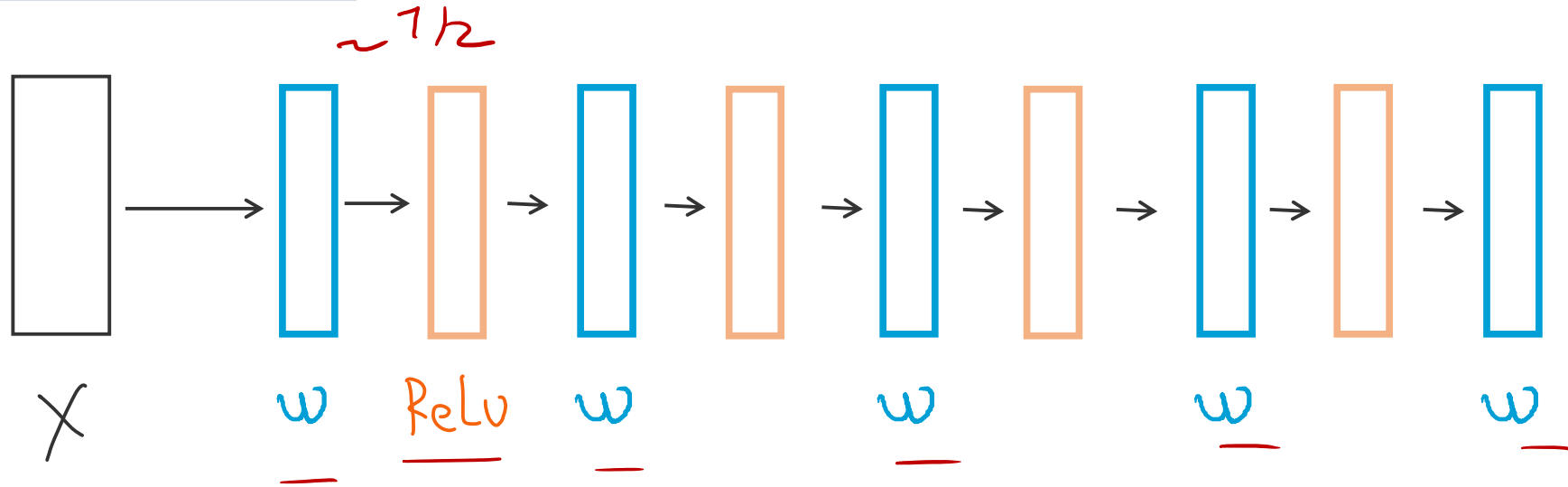
```
a = 1 / sqrt(in_num)  
W = a * random(in_num, out_num)
```



```
torch.nn.init.xavier_normal_(
```

He initialization

```
a = 1 / sqrt(in_num / 2)  
W = a * random(in_num, out_num)
```



```
torch.nn.init.kaiming_normal_(
```

Обновление параметров

Update rules

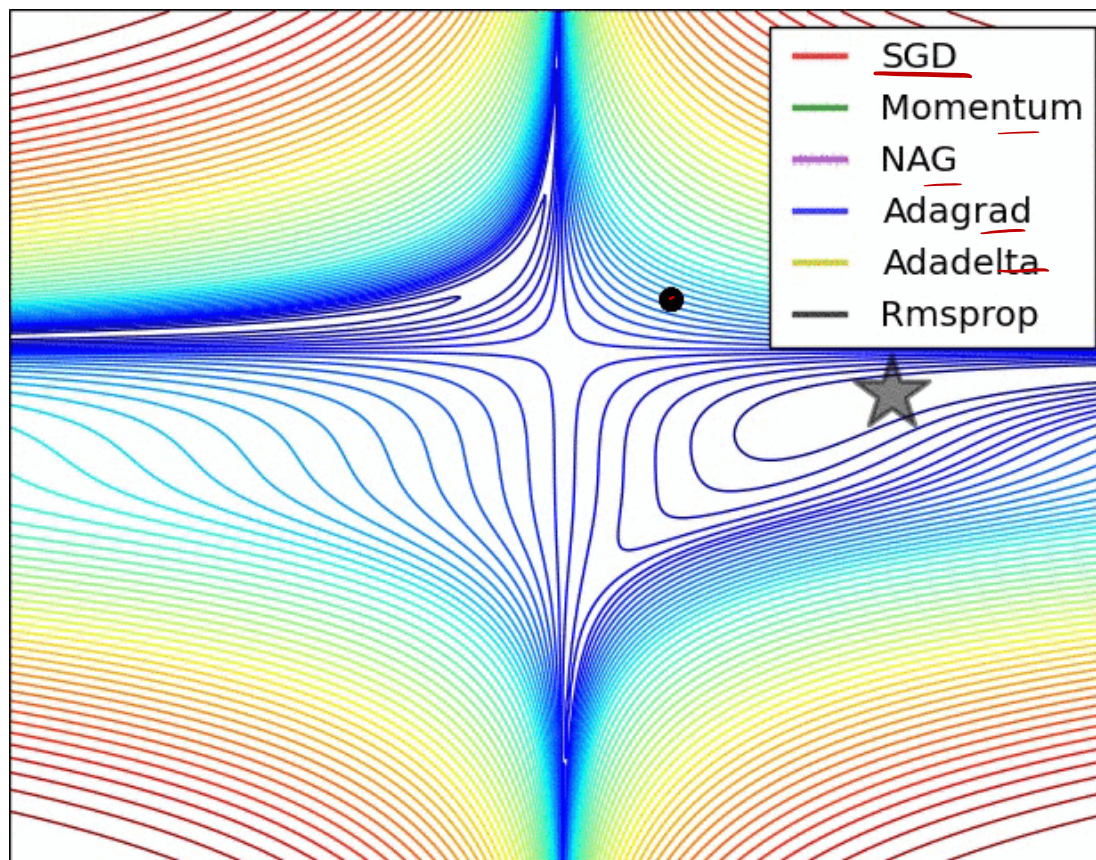
$$\underline{w} = w - \underline{\text{learning_rate}} * \text{gradient}$$

Обновить параметры

$$\underline{\vec{w}^1} = \vec{w}^1 - \eta \vec{\nabla}_{w^1} L$$

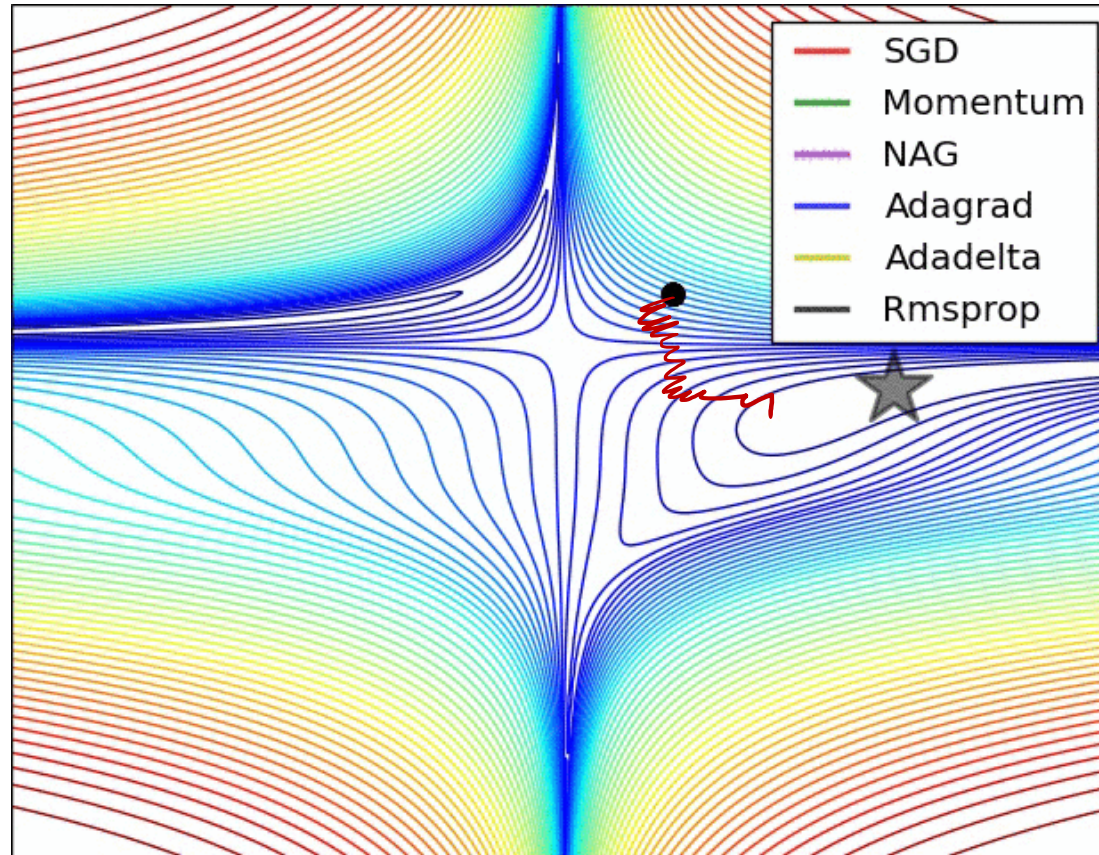
...

$$\underline{\vec{w}^n} = \vec{w}^n - \eta \vec{\nabla}_{w^n} L$$

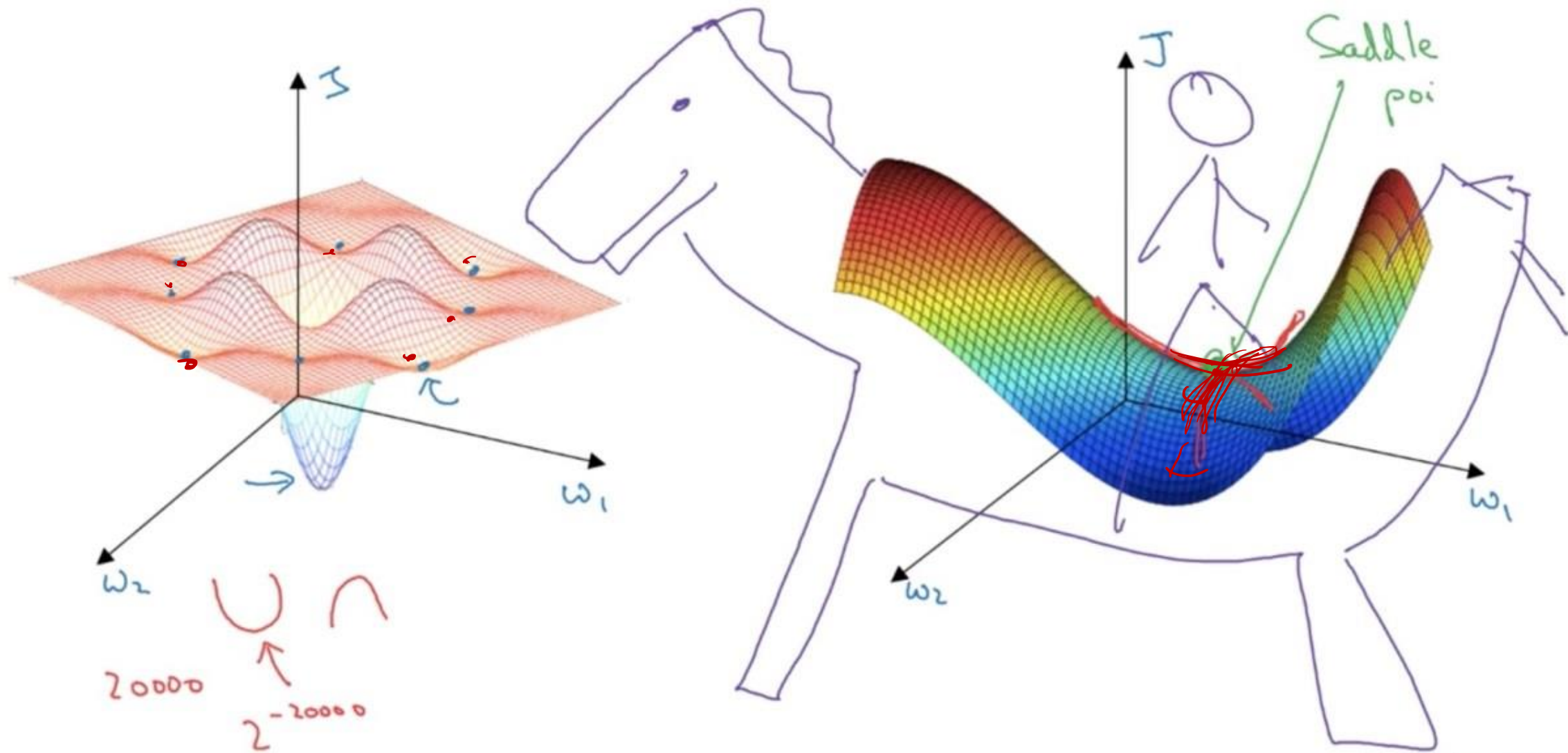


Momentum *(импульс)*

0.9-0.99
velocity = momentum * velocity - learning_rate * gradient,
w = w + velocity



Local optima in neural networks

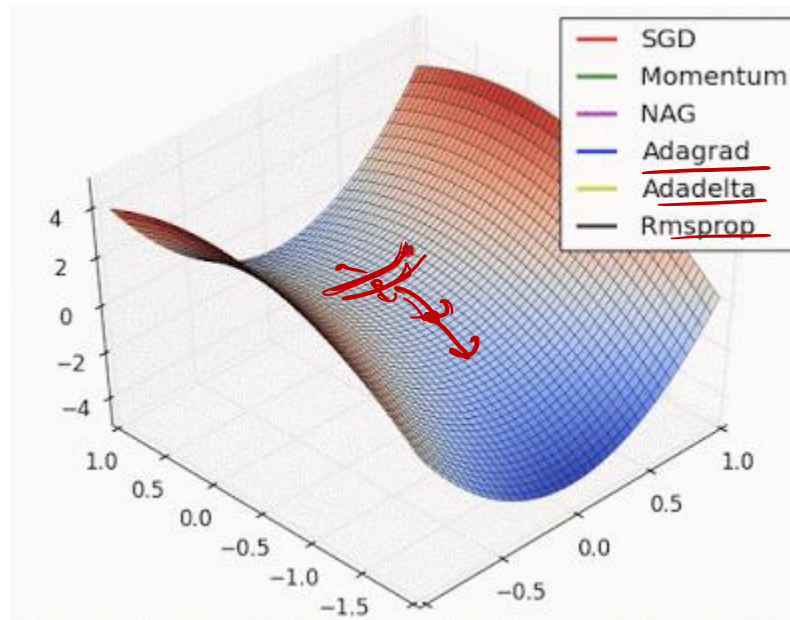


Andrew Ng

Adagrad

```
accumulated += |gradient| ** 2  
adaptive_learning_rate = learning_rate / sqrt(accumulated)
```

```
w := w - adaptive_learning_rate * gradient
```



```
torch.optim.Adagrad()
```


RMSProp

```
torch.optim.RMSprop(
```

0.9 - 0.99

```
accumulated = rho * accumulated + (1-rho) * gradient ** 2  
adaptive_learning_rate = learning_rate / sqrt(accumulated)
```

```
w = w - adaptive_learning_rate * gradient
```

[Neural Networks for Machine Learning, Lecture 6.5 – rmsprop'12](#)

Adam

```
torch.optim.Adam(
```

```
velocity = beta1 * velocity + (1-beta1) * gradient  
accumulated = beta2 * accumulated + (1-beta2) * gradient ** 2
```

```
adaptive_learning_rate = learning_rate / sqrt(accumulated)
```

```
w = w - adaptive_learning_rate * velocity
```

*RMS Prop
+ Momentum
Peko metgyen*

[Adam: A Method for Stochastic Optimization'14](#)

В следующий раз

