

# BACKPROPAGATION ADVANCED

## Постановка задачи

### Дано:

$$\mathcal{X} = (X_1, \dots, X_k)$$

$$\mathcal{A} = (A_1, \dots, A_k)$$

$$(\mathcal{X}, \mathcal{A})$$

$$W$$

$$N(W, X)$$

$$Y = N(W, X)$$

$$D(Y, A) = \sum_{j=1}^m (Y[j] - A[j])^2$$

$$D_i(Y) = D(Y, A_i)$$

$$E_i(W) = D_i(N(W, X_i))$$

$$E(W) = \sum_{i=1}^k E_i(W)$$

### Найти:

вектор  $W$  такой, что  $E(W) \rightarrow \min$  (обучение на всей выборке)

вектор  $W$  такой, что  $E_i(W) \rightarrow \min$  (обучение на одном примере)

входные вектора,  $X_i \in \mathbb{R}^n$

правильные выходные вектора,  $A_i \in \mathbb{R}^m$

обучающая выборка

вектор весов нейронной сети

функция, соответствующая нейронной сети

ответ нейронной сети,  $Y \in \mathbb{R}^m$

функция ошибки

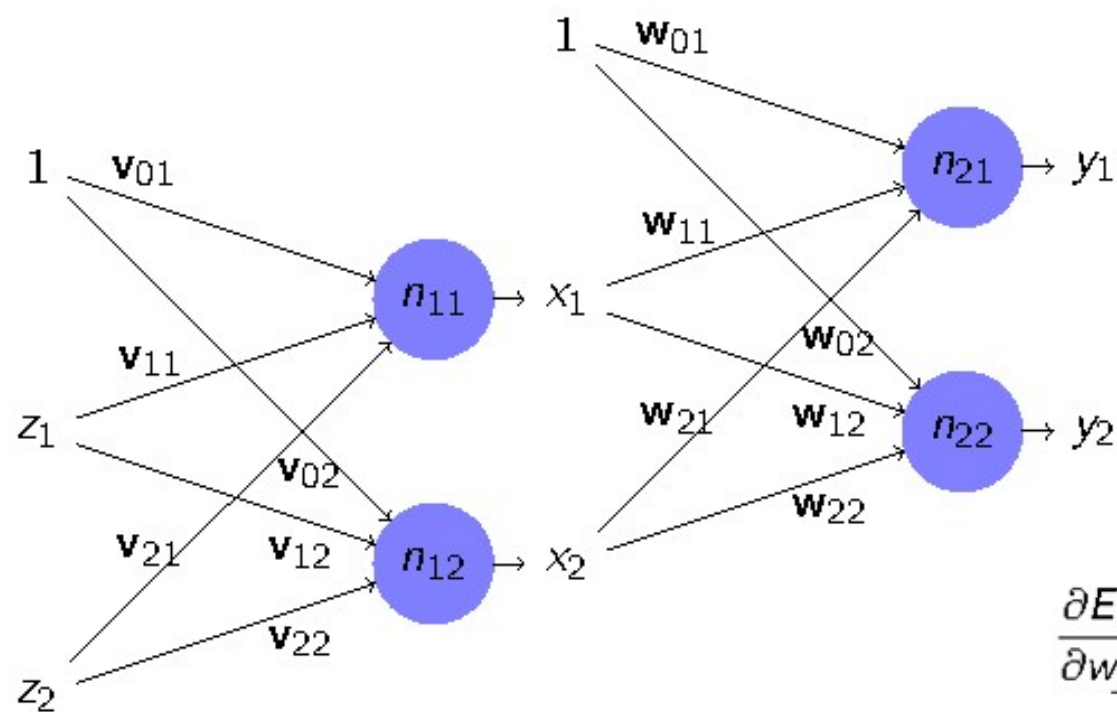
функция ошибки на  $i$ -ом примере

ошибка сети на  $i$ -ом примере

ошибка сети на всей обучающей выборке

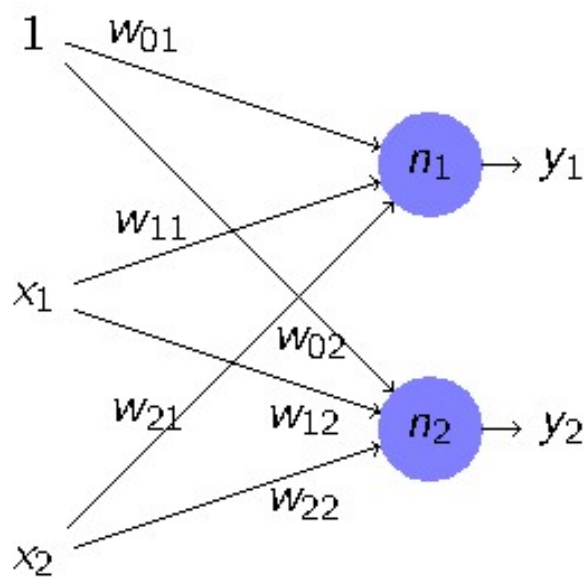
$$\begin{aligned}
\frac{\partial E(\vec{w})}{\partial w_j} &= \frac{\partial}{\partial w_j} \sum_i (h_{\vec{w}}(\vec{x}_i) - y_i)^2 \\
&= \sum_i 2(h_{\vec{w}}(\vec{x}_i) - y_i) \frac{\partial}{\partial w_j} (h_{\vec{w}}(\vec{x}_i) - y_i) \\
&= \sum_i 2(h_{\vec{w}}(\vec{x}_i) - y_i) \frac{\partial}{\partial w_j} \sigma(\vec{x}_i \cdot \vec{w}) \\
&= \sum_i 2(h_{\vec{w}}(\vec{x}_i) - y_i) \sigma'(\vec{x}_i \cdot \vec{w}) \frac{d}{dw_j} \vec{x}_i \cdot \vec{w} \\
&= \sum_i 2(h_{\vec{w}}(\vec{x}_i) - y_i) \sigma'(\vec{x}_i \cdot \vec{w}) \frac{d}{dw_j} \sum_{k=1}^n x_{i,k} w_k \\
&= 2 \sum_i (h_{\vec{w}}(\vec{x}_i) - y_i) \sigma'(\vec{x}_i \cdot \vec{w}) x_{i,j}
\end{aligned}$$

## Обратное распространение ошибки



$$\frac{\partial E_k}{\partial w_{ji}} = \sum_{l=1}^n \frac{\partial D_k}{\partial y_l} \frac{\partial y_l}{\partial w_{ji}} = 2(y_i - a_i) f'(S_i) x_j$$

## Обратное распространение ошибки



$$D_k(y_1, \dots, y_n) = (y_i - a_i)^2 + \dots + (y_n - a_n)^2$$

$$\frac{\partial D_k}{\partial y_i} = 2(y_i - a_i)$$

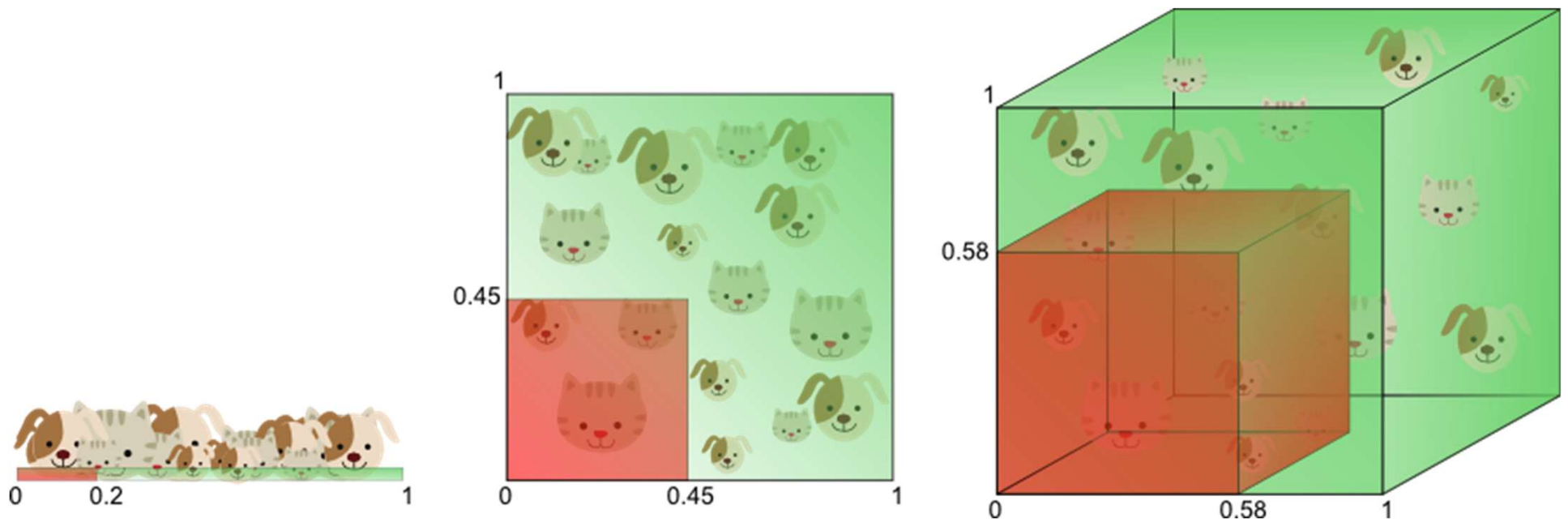
$$S_i = \sum_{j=0}^m x_j w_{ji} \quad y_i = f(S_i) \quad \frac{\partial y_i}{\partial x_j} = f'(S_i) w_{ji}$$

$$\begin{aligned} \frac{\partial D_k}{\partial x_j} &= \sum_{i=1}^n \frac{\partial D_k}{\partial y_i} \frac{\partial y_i}{\partial x_j} = \\ &= 2 \sum_{i=1}^n (y_i - a_i) f'(S_i) w_{ji} \end{aligned}$$

# Algorithm of BP

```
initialize network weights (often small random values)
do
  forEach training example named ex
    prediction = neural-net-output(network, ex) // forward pass
    actual = teacher-output(ex)
    compute error (prediction - actual) at the output units
    compute  $\Delta w_h$  for all weights from hidden layer to output layer // backward pass
    compute  $\Delta w_i$  for all weights from input layer to hidden layer // backward pass continued
    update network weights // input layer not modified by error estimate
  until all examples classified correctly or another stopping criterion satisfied
return the network
```

# The curse of dimensionality



# Plan

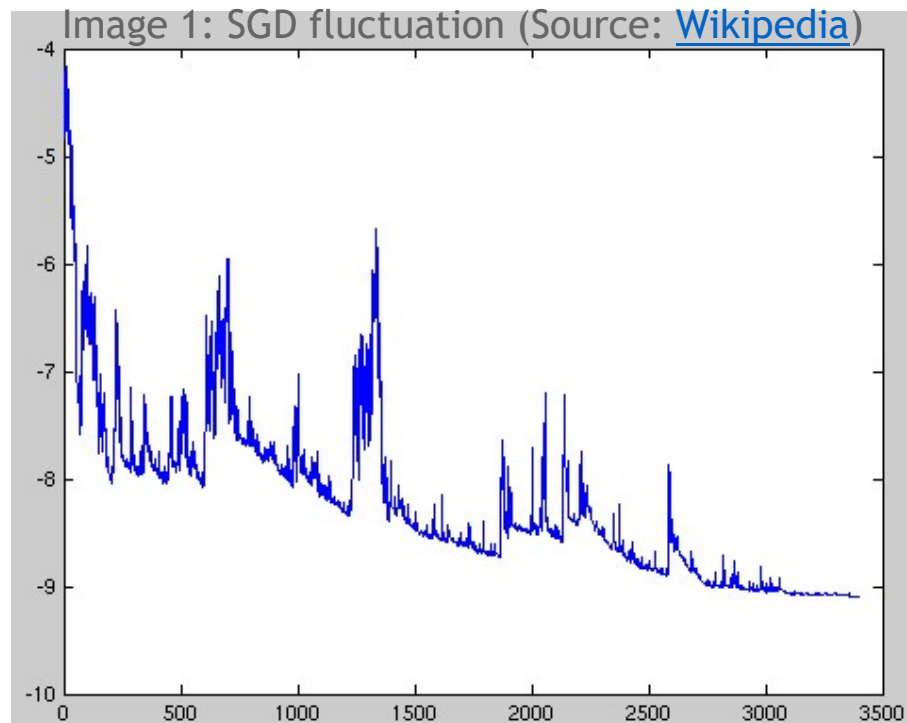
- Gradient descent variants
- Challenges
- Gradient descent optimization algorithms
- Parallelizing and distributing SGD
- Additional strategies for optimizing SGD



# Batch gradient descent

Stochastic:  $\theta = \theta - \eta \cdot \nabla_{\theta} J(\theta).$

Batch:  $\theta = \theta - \eta \cdot \nabla_{\theta} J(\theta; x^{(i)}; y^{(i)}).$



# Batch vs Stochastic

```
for i in range(nb_epochs):  
    params_grad = evaluate_gradient(loss_function, data, params)  
    params = params - learning_rate * params_grad
```

```
for i in range(nb_epochs):  
    np.random.shuffle(data)  
    for example in data:  
        params_grad = evaluate_gradient(loss_function, example, params)  
        params = params - learning_rate * params_grad
```

## Challenges: Vanilla mini-batch gradient descent

- Choosing a proper learning rate can be difficult
- Learning rate schedules try to adjust the learning rate during.
- Additionally, the same learning rate applies to all parameter updates.
- Suboptimal local minima.

# Mini-batch gradient descent

$$\theta = \theta - \eta \cdot \nabla_{\theta} J(\theta; x^{(i:i+n)}; y^{(i:i+n)}).$$

```
for i in range(nb_epochs):  
    np.random.shuffle(data)  
    for batch in get_batches(data, batch_size=50):  
        params_grad = evaluate_gradient(loss_function, batch, params)  
        params = params - learning_rate * params_grad
```

# Gradient descent optimization algorithms

Momentum

Nesterov accelerated gradient

Adagrad

Adadelat

RMSprop

Adam

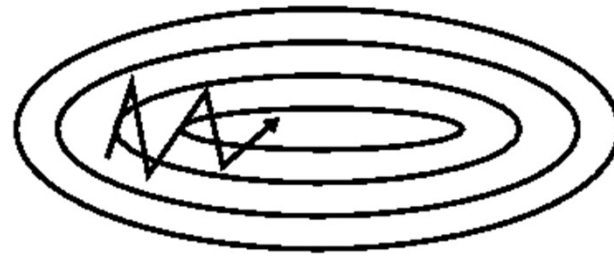
AdaMax

Nadam

Visualization of algorithms

Which optimizer to choose?

# Momentum



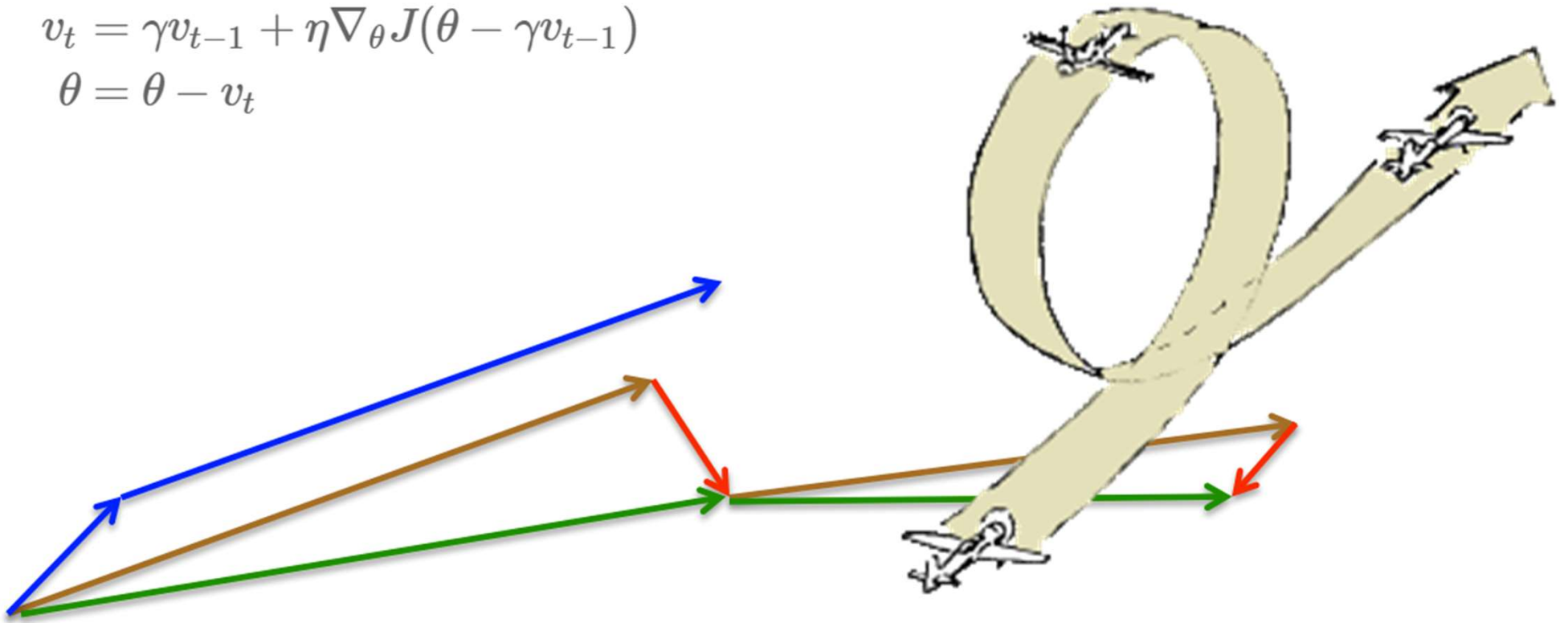
$$v_t = \gamma v_{t-1} + \eta \nabla_{\theta} J(\theta)$$

$$\theta = \theta - v_t$$

# Nesterov accelerated gradient

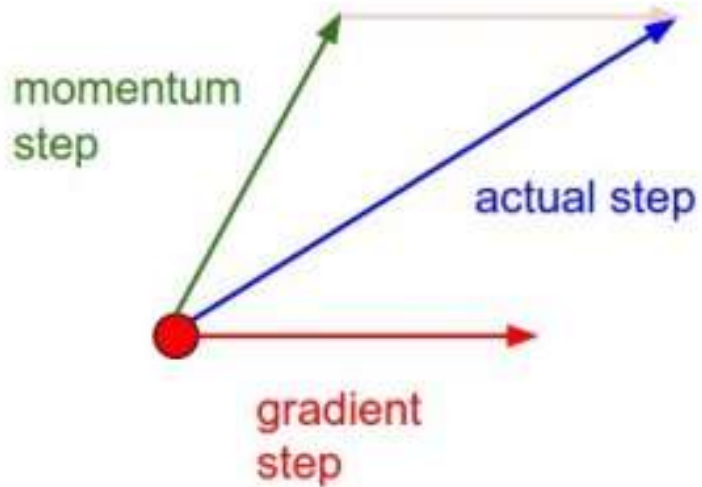
$$v_t = \gamma v_{t-1} + \eta \nabla_{\theta} J(\theta - \gamma v_{t-1})$$

$$\theta = \theta - v_t$$

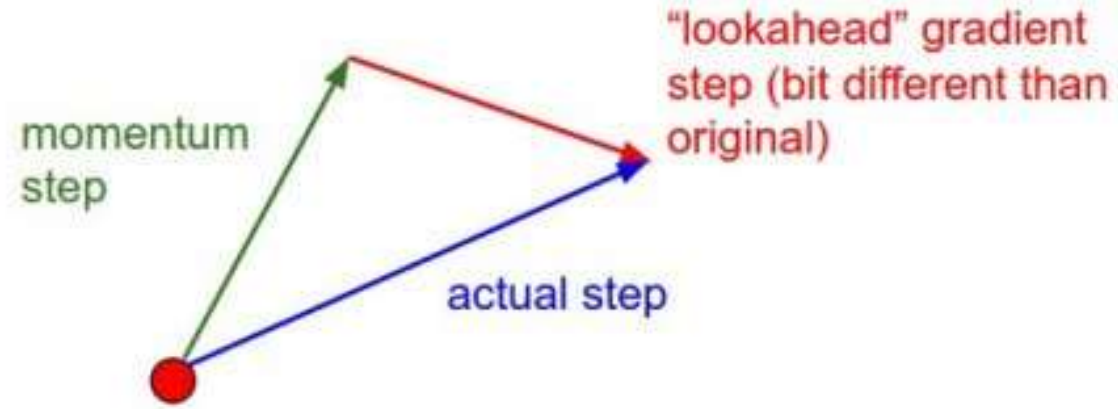


# Momentum vs Nesterov Acc. Grad.

Momentum update



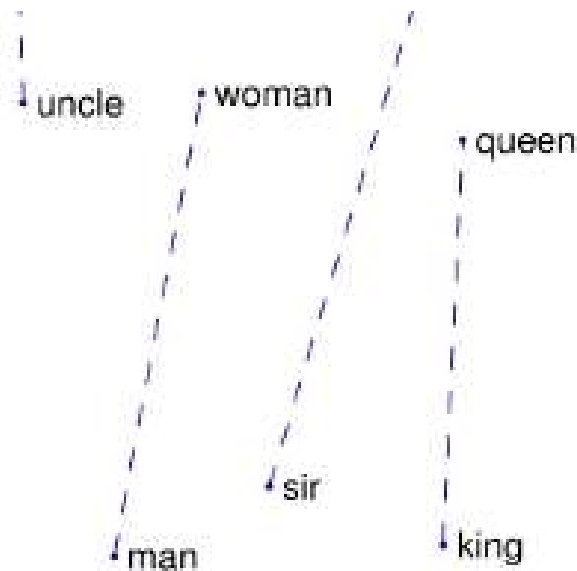
Nesterov momentum update





# Adagrad

It adapts the learning rate to the parameters, performing larger updates for infrequent and smaller updates for frequent parameters.



# GOOGLE'S ARTIFICIAL BRAIN LEARNS TO FIND CAT VIDEOS



# Adagrad

## Adaptive Gradients

$g_{t,i} = \nabla_{\theta} J(\theta_i)$ . the gradient of the objective function w.r.t. to the parameter  $\theta$  at time step  $t$

$\theta_{t+1,i} = \theta_{t,i} - \eta \cdot g_{t,i}$ . The SGD update for every parameter  $\theta$  at each time step  $t$  then becomes



$$\theta_{t+1,i} = \theta_{t,i} - \frac{\eta}{\sqrt{G_{t,ii} + \epsilon}} \cdot g_{t,i}.$$



# Adagrad

$G_t \in \mathbb{R}^{d \times d}$  diagonal matrix the sum of the squares of the gradients

$$G = \sum_{\tau=1}^t g_{\tau} g_{\tau}^{\top}$$

$$\theta_{t+1,i} = \theta_{t,i} - \frac{\eta}{\sqrt{G_{t,ii}} + \epsilon} \cdot g_{t,i} \quad \longrightarrow \quad \theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{G_t} + \epsilon} \odot g_t.$$

# Adadelata

Adadelata is an extension of Adagrad that seeks to reduce its aggressive, monotonically decreasing learning rate. Instead of accumulating all past squared gradients, Adadelata restricts the window of accumulated past gradients to some fixed size  $w$ .

# Adadelta

$$E[g^2]_t = \gamma E[g^2]_{t-1} + (1 - \gamma)g_t^2.$$

$$\Delta\theta_t = -\eta \cdot g_{t,i}$$

$$\theta_{t+1} = \theta_t + \Delta\theta_t$$

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{G_t + \epsilon}} \odot g_t. \quad \longrightarrow \quad \Delta\theta_t = -\frac{\eta}{\sqrt{E[g^2]_t + \epsilon}} g_t.$$

# Adadelta

$$\Delta\theta_t = -\frac{\eta}{RMS[g]_t} g_t.$$

$$E[\Delta\theta^2]_t = \gamma E[\Delta\theta^2]_{t-1} + (1 - \gamma) \Delta\theta_t^2. \quad (*)$$

$$RMS[\Delta\theta]_t = \sqrt{E[\Delta\theta^2]_t + \epsilon}.$$

$$\Delta\theta_t = -\frac{RMS[\Delta\theta]_{t-1}}{RMS[g]_t} g_t$$

$$\theta_{t+1} = \theta_t + \Delta\theta_t$$

# RMSprop

$$E[g^2]_t = 0.9E[g^2]_{t-1} + 0.1g_t^2$$
$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{E[g^2]_t + \epsilon}} g_t$$



Geoff Hinton

# Adam

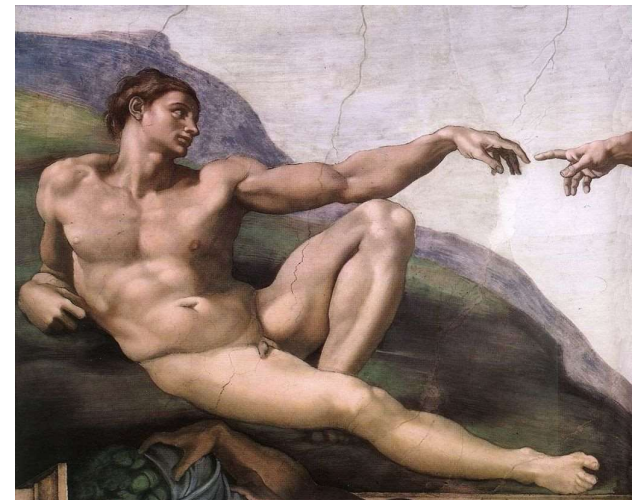
## Adaptive Moment Estimation

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$$
$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$$

The Analogue between Adam and (RMSProp, Ada\*)

Init m and v: 0

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}$$
$$\hat{v}_t = \frac{v_t}{1 - \beta_2^t}$$
$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \hat{m}_t.$$





# Nadam

## Nesterov-accelerated Adaptive Moment Estimation

$$g_t = \nabla_{\theta_t} J(\theta_t)$$

$$m_t = \gamma m_{t-1} + \eta g_t$$

$$\theta_{t+1} = \theta_t - m_t$$

$$\theta_{t+1} = \theta_t - (\gamma m_{t-1} + \eta g_t)$$

Recall of momentum

Recall of NAG

$$g_t = \nabla_{\theta_t} J(\theta_t - \gamma m_{t-1})$$

$$m_t = \gamma m_{t-1} + \eta g_t$$

$$\theta_{t+1} = \theta_t - m_t$$

# Nadam

$$g_t = \nabla_{\theta_t} J(\theta_t)$$

$$m_t = \gamma m_{t-1} + \eta g_t$$

Another way of NAG

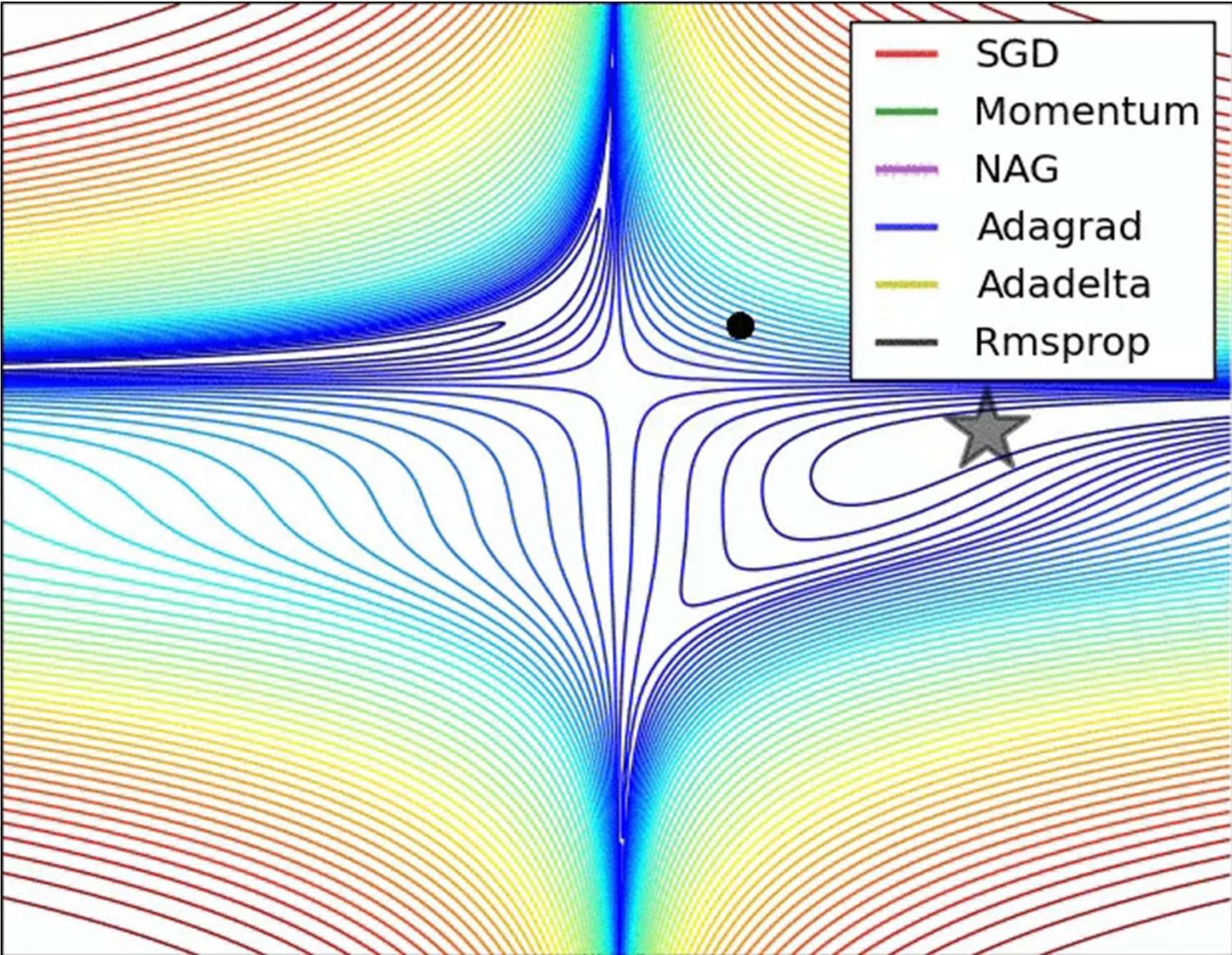
$$\theta_{t+1} = \theta_t - (\gamma m_t + \eta g_t)$$

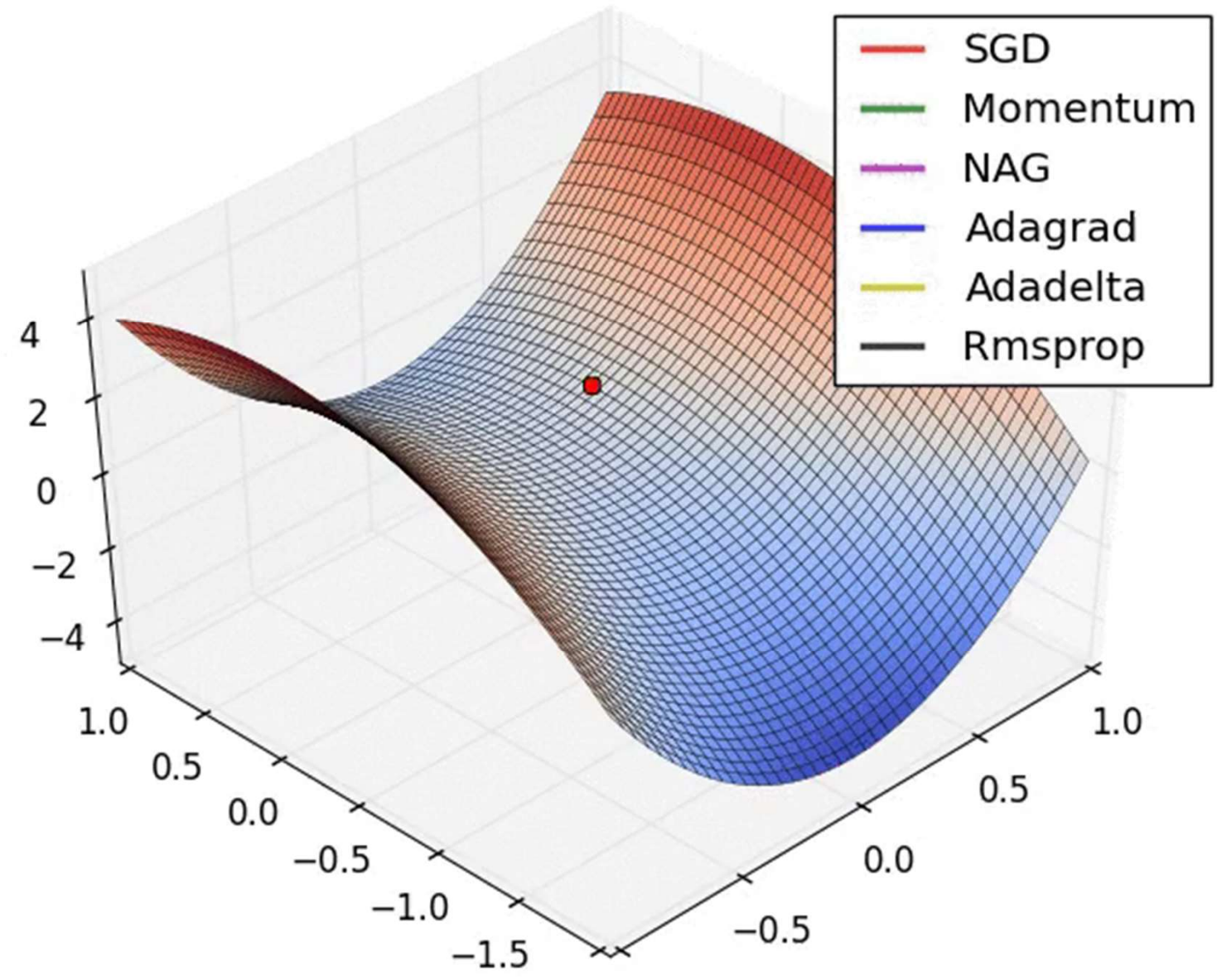
$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$$

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}$$

Nadam

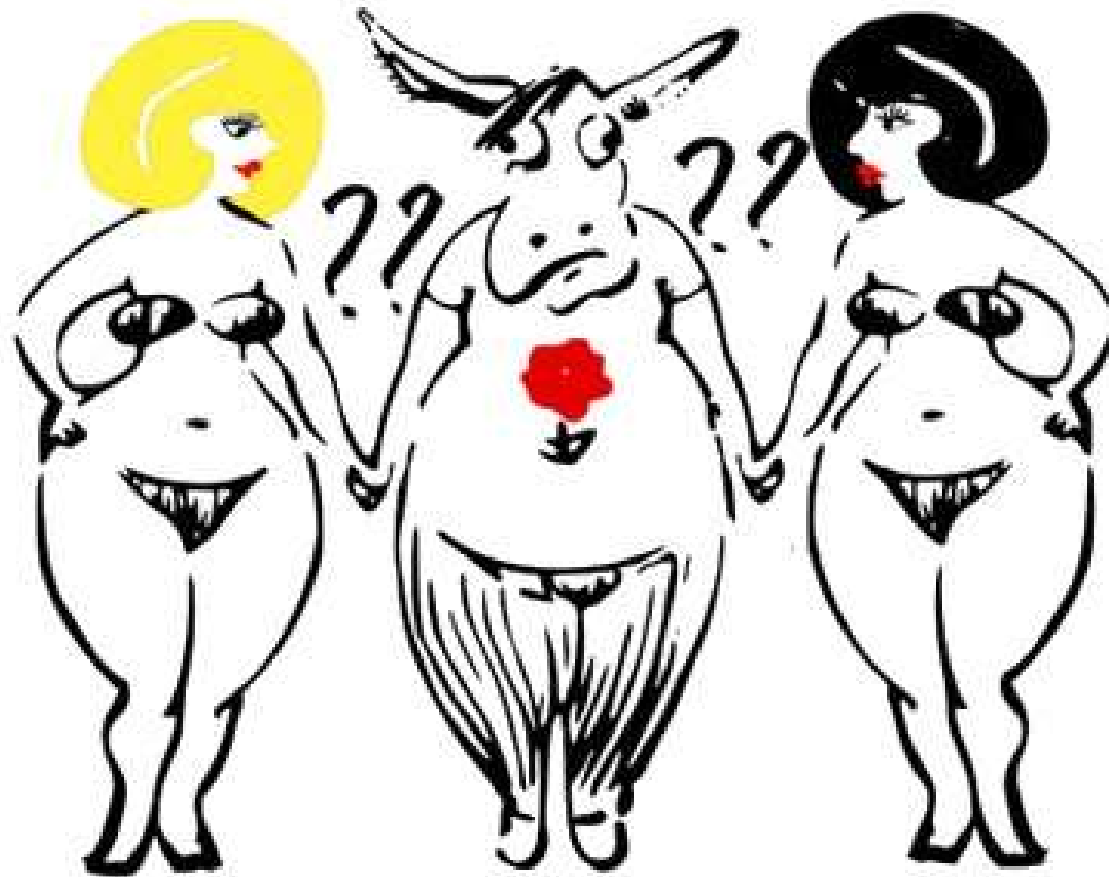
$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \hat{m}_t$$





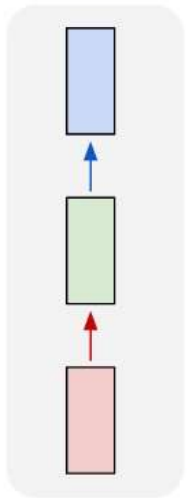


Which one to use?

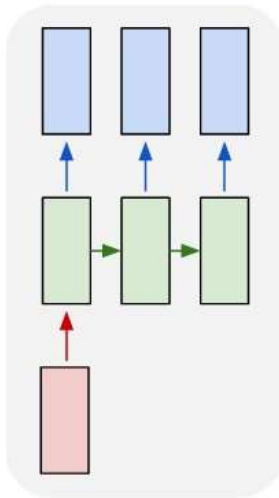


# RNN

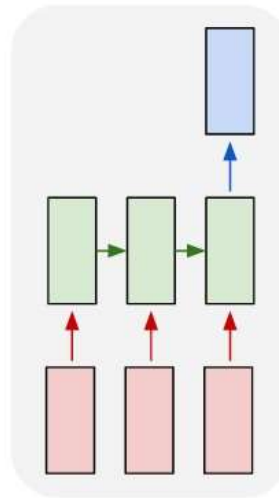
one to one



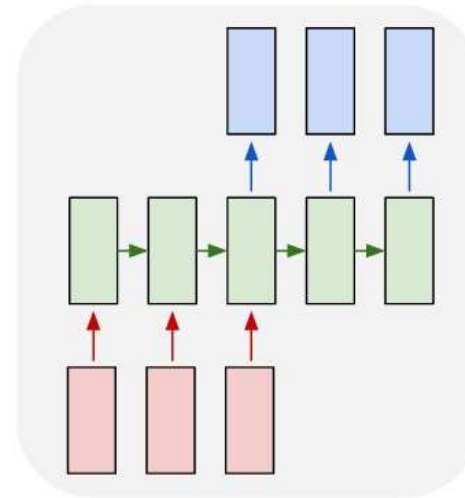
one to many



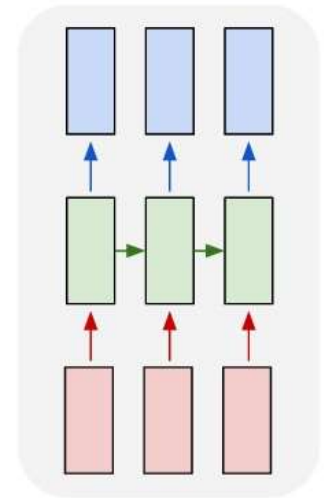
many to one



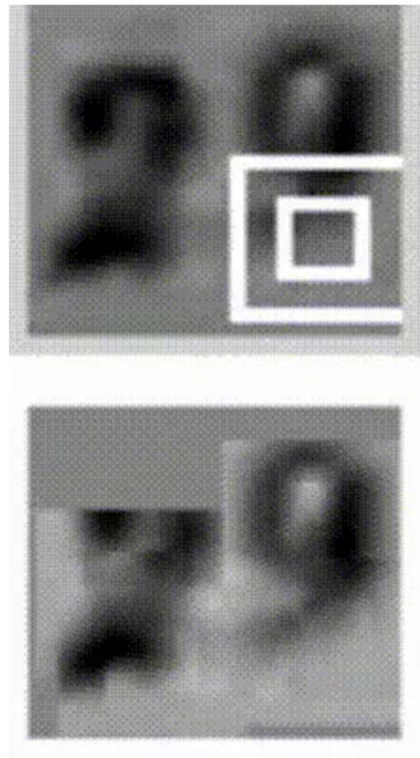
many to many



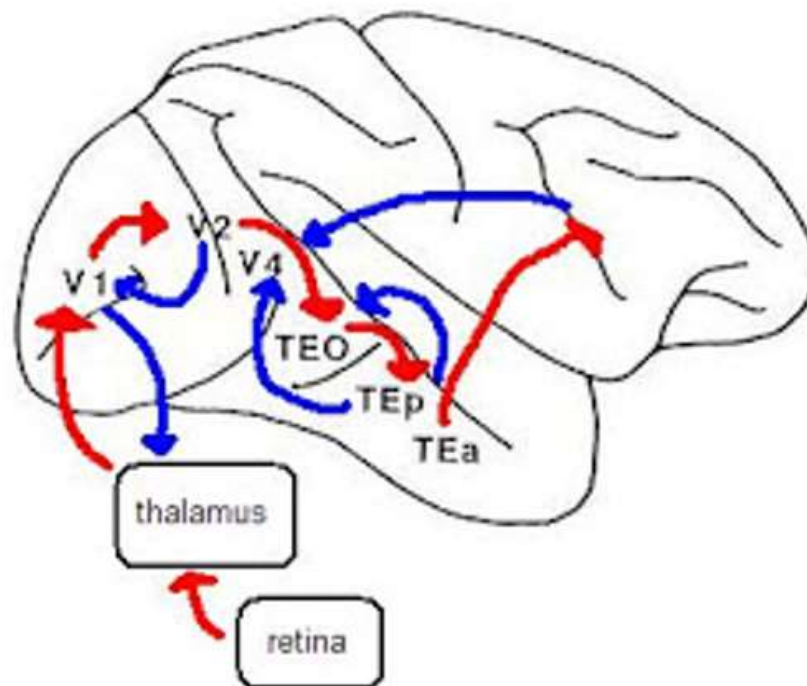
many to many



# Example



# Временные ряды: Рабочая память



Long Short-Term Memory  
(LSTM)

Recurrent Neural Net

Recurrent Processing - concious?

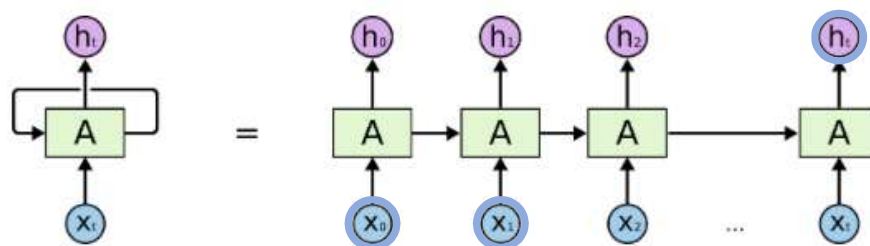


# RNN Code Example

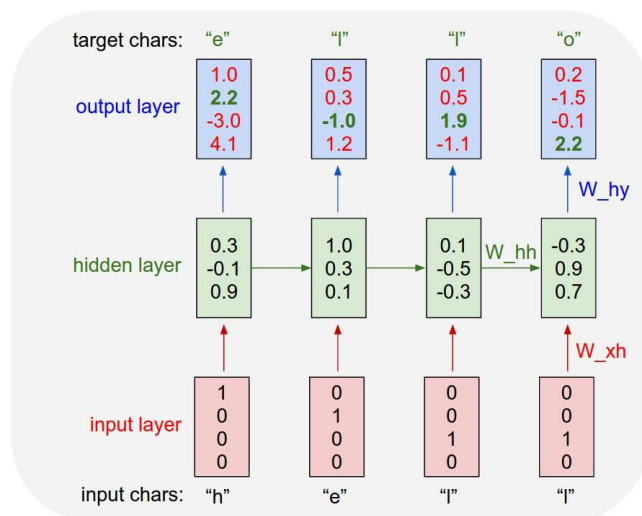
```
rnn = RNN()  
y = rnn.step(x) # x is an input vector, y is the RNN's output vector
```

```
class RNN: # ...  
    def step(self, x):  
        # update the hidden state  
        self.h = np.tanh(np.dot(self.W_hh, self.h) + np.dot(self.W_xh, x))  
        # compute the output vector  
        y = np.dot(self.W_hy, self.h)  
        return y
```

# Рекуррентные сети (Long Short-Term Memory)

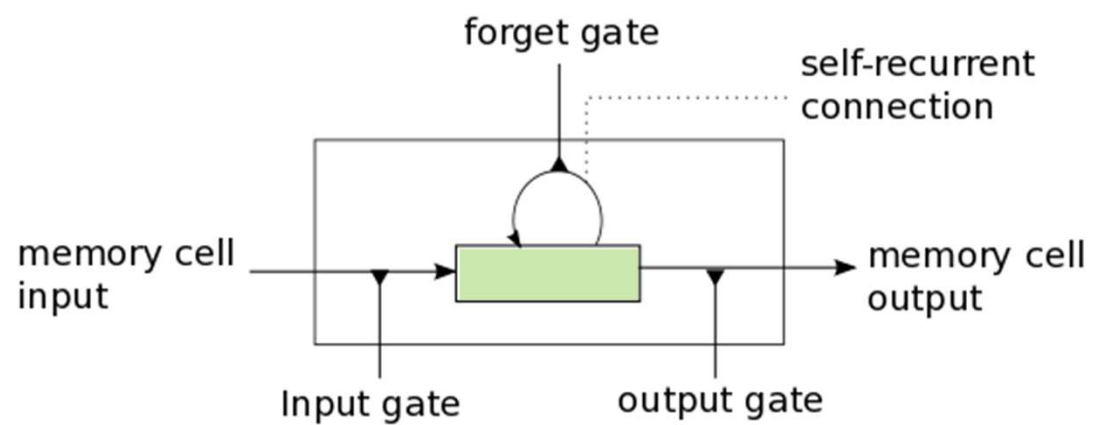


Временная  
инвариантность



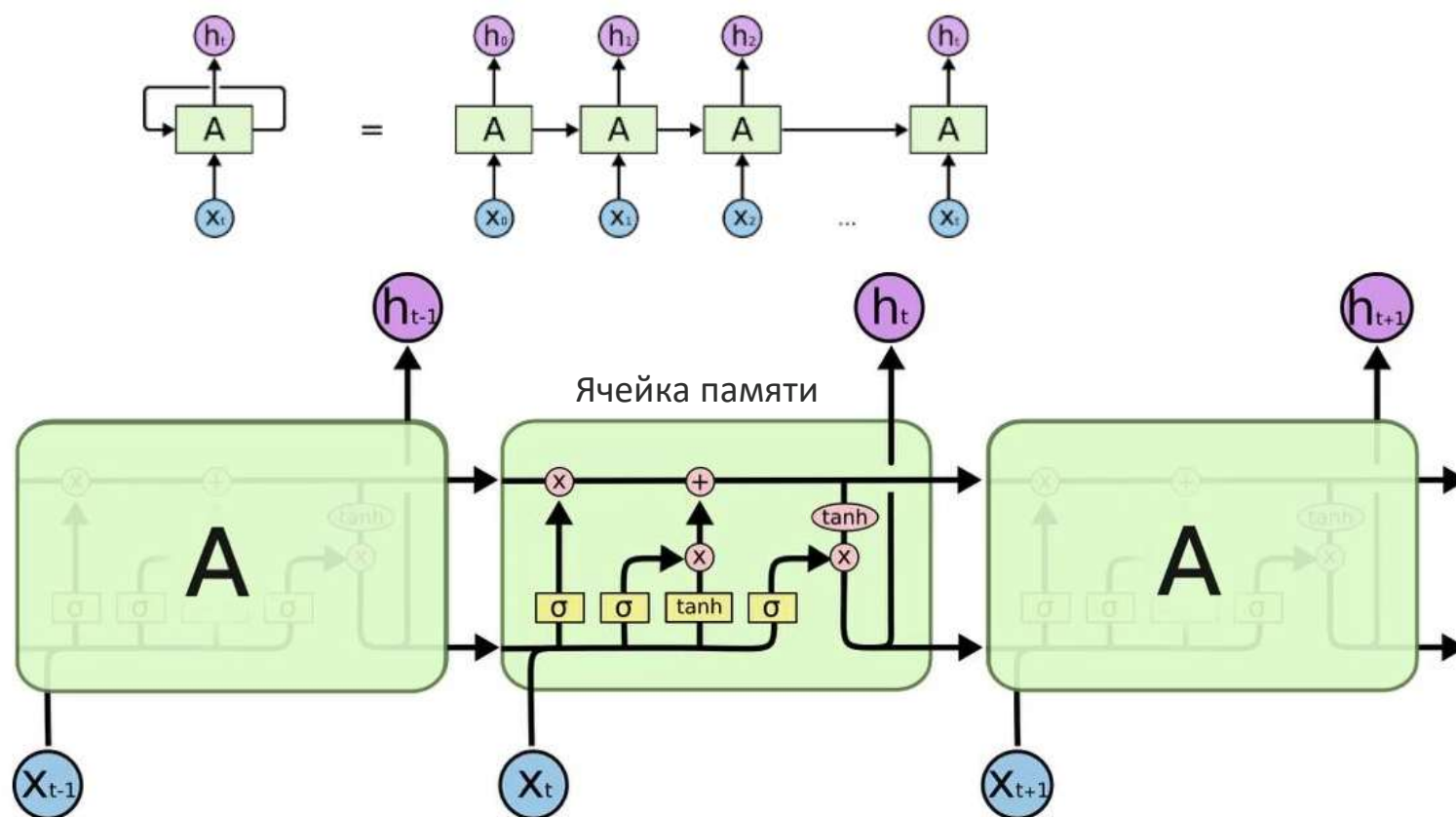
Hochreiter (1997)  
*Long Short-Term Memory*

# LSTM Architecture



Ячейка памяти

# Рекуррентные сети (Long Short-Term Memory)



# Shakespeare\Wiki\Latex toy models

tyntd-iafhatawiao hr demot lytdws e ,tfti, astai f ogoh eoase  
rrranbyne 'nhthnee e plia tk lrgd t o idoe ns, smtt h ne etie  
h, hregtrs nigtike, aoaenns lng

"Tmont thithey" fomesscerliund Keushey. Thom here sheulke,  
anmerenith ol sivh I lalterthend Bleipile shuwy fil on aseterlome  
coaniogennc Phe lism thond hon at. MeiDimorotion in ther thize."

we counter. He stutn co des. His stanted out one ofler that  
concoctions and was to gearang reay Jotrets and with fre colt off paitt  
thin wall. Which das stimn

# LSTM Vizualization

t t p : / / w w w . y n e t n e w s . c o m / ] E n g l i s h - l a n g u a g e w e b s i t e o f I s r a e l ' s l a r  
t p : / / w w w b a c a h e t s . c o m / - x g l i s h l i n g u a g e s a i r s i t e o f t s l a e l i s s i n g  
d : x n e . w a e a . a w a t o a . s & n t i a c a - s a r d e e l h o a n t b i s a n f a n r e i f ' a a t d  
m w - 2 p i i i s o e s s i s . / e r n . c ] ( d c e e n e p e s a a i k i i e e l e d h , i r t h r a o n s e , c o s e  
d r . < : a h b - n p t w t . x i g h / m a ) T v d r y z i c o u e d l s u : t h a - o o t u , s t u i f l v e p e r y  
s t p , t c o a 2 d r u l w o c l e n s r ] p . l l v a o d , , e y t c - n d m - o i b u v s ] b b i m s u l t a t t l y b n

g e s t n e w s p a p e r ' ' [ [ Y e d i o t h A h r o n o t h ] ] ' ' ' ' H e b r e w - l a n g u a g e p e r i o d  
e t a a w s p a p e r s o [ [ T e l t i ( f e a n e m t i ) ' ' \* ' ' [ e r r e w s l e n g u a g e : a r o s o d i  
i r s c o e e n a i T T h A o a i n n h S r m u w ] e y s [ ' i n e i a ' s i w d d e ' h s o l r i f r :  
u s . . s e t l g o r s . a s a t C a r e e g ' a C l r i s z ] i e ' : : , # : T A a a a a t B a s e e i l o ' i a n f v l  
- t u a e v r t i d , t B A m S u s y u t ] ] A s a o i g s ] ] , . : s M B o l o u s : T o u a - n : d w o a p n u  
a , d , i i u i t i c p . ] ( l S v H v t u s u i e D n o e g a n o . , ] : { C C u i b o h e C y b k s l s : r - e p c n t s

i c a l s : ' ' ' \* ' ' [ [ G l o b e s ] ] ' ' [ h t t p : / / w w w . g l o b e s . c o . i l / ] b u s i n e s s d a  
c a l : ' ' ' \* ' ' [ T a a b a ] ' ' ( [ t t p : / / w w w . b u o b a l . c o m u n / s A - y t i n e s s a e t  
s t l ' [ h A e o v e l t s a h a d : x g e . w a o i r . r t o a . e l . i T & a i e g e o o y  
t t ' ' ' & [ & & m C o e r o n e ' : : , i ' o d w . , : n i i i s a a u e . e n i / o m l c C . ( e f t g i r i i u  
a ' n : , C : & : # \* : a f D r u s u ] l , . o m e l p < , d h a ; d e u o o t / i h n c s i f S , u r h o s t , t u n  
n k i < ] : & 1 1 s T G u i t r s i , : b a c m r - x t p o b - g r e s i s l e r l n a f a D ] l o s p t a d , i f r m

i l y \* ' ' [ [ H a a r e t z ] H a ' A r e t z ] ] ' ' [ h t t p : / / w w w . h a a r e t z . c o . i l / ] R e l a t i v  
l y \* ' ' [ [ T e r r d n F e r a n t a h ] ] ' ' ( [ t t p : / / w w w . b o n m d s t . c o m u n / s - e s a t e o i  
r e ' ' ' h A i l n n t t e H a l s r c n o l ' s a h a d : x n e . w a a m r t d h e o h . o l . c & o p i n i v e  
k i . : \* s C O S a n l t h i T i m ' l i ] e : , i m c d w - 2 p h i i s e r d i t . i n a / c m f i . ( a f l c a n a  
d s - ! [ t B T C o m m g d ] ] W o n a a e , : . b a e r r . < t a i b - d u l c n n c / a r n e s i ] l i c e y s t o  
n d s # & : G l D u v c c s a o S u c l t e l ] z | , : o ' o m t ] , : e o a 2 n i v f s r o o e i u n a l a ) u v v r o

# LSTM vs GRU

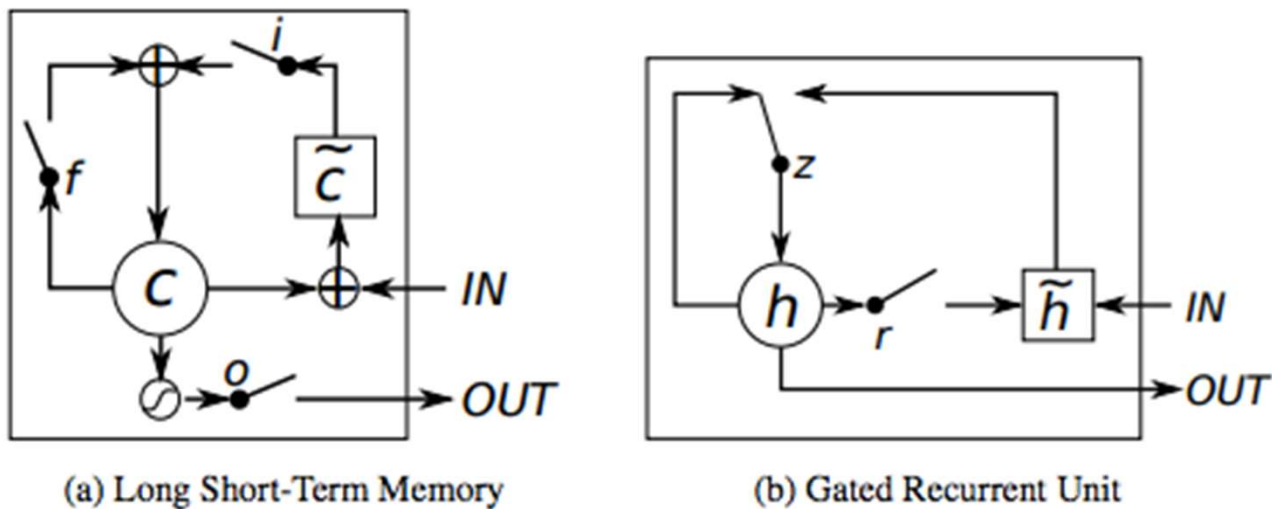


Figure 1: Illustration of (a) LSTM and (b) gated recurrent units. (a)  $i$ ,  $f$  and  $o$  are the input, forget and output gates, respectively.  $c$  and  $\tilde{c}$  denote the memory cell and the new memory cell content. (b)  $r$  and  $z$  are the reset and update gates, and  $h$  and  $\tilde{h}$  are the activation and the candidate activation.

# LSTM output gate

Unlike to the recurrent unit which simply computes a weighted sum of the input signal and applies a nonlinear function, each  $j$ -th LSTM unit maintains a memory  $c_t^j$  at time  $t$ . The output  $h_t^j$ , or the activation, of the LSTM unit is then

$$h_t^j = o_t^j \tanh \left( c_t^j \right),$$

where  $o_t^j$  is an *output gate* that modulates the amount of memory content exposure. The output gate is computed by

$$o_t^j = \sigma \left( W_o \mathbf{x}_t + U_o \mathbf{h}_{t-1} + V_o \mathbf{c}_t \right)^j,$$

where  $\sigma$  is a logistic sigmoid function.  $V_o$  is a diagonal matrix.



# LSTM memory gate

The memory cell  $c_t^j$  is updated by partially forgetting the existing memory and adding a new memory content  $\tilde{c}_t^j$  :

$$c_t^j = f_t^j c_{t-1}^j + i_t^j \tilde{c}_t^j, \quad (4)$$

where the new memory content is

$$\tilde{c}_t^j = \tanh (W_c \mathbf{x}_t + U_c \mathbf{h}_{t-1})^j .$$

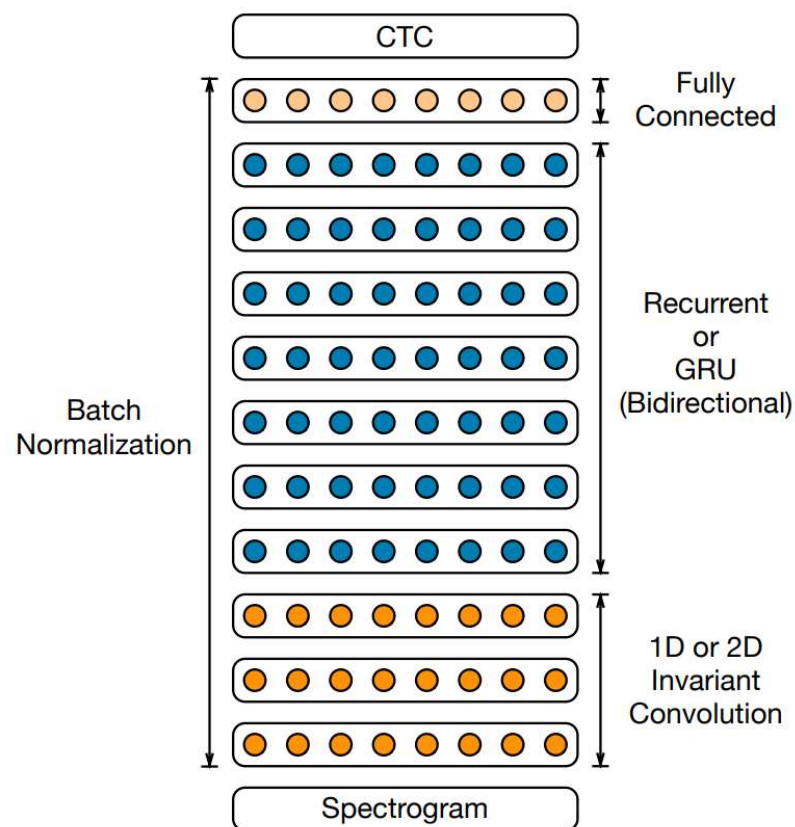
## LSTM Input gate

The extent to which the existing memory is forgotten is modulated by a *forget gate*  $f_t^j$ , and the degree to which the new memory content is added to the memory cell is modulated by an *input gate*  $i_t^j$ . Gates are computed by

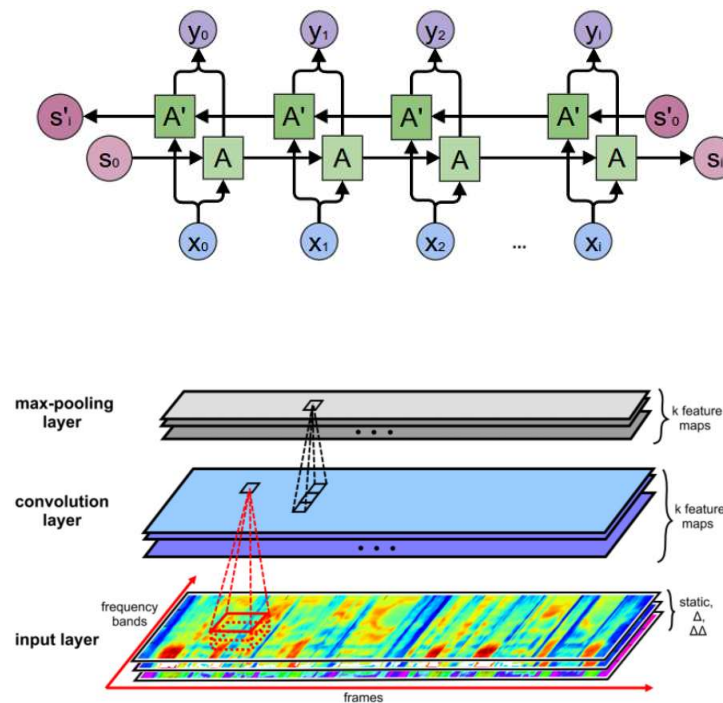
$$\begin{aligned} f_t^j &= \sigma (W_f \mathbf{x}_t + U_f \mathbf{h}_{t-1} + V_f \mathbf{c}_{t-1})^j, \\ i_t^j &= \sigma (W_i \mathbf{x}_t + U_i \mathbf{h}_{t-1} + V_i \mathbf{c}_{t-1})^j. \end{aligned}$$

Note that  $V_f$  and  $V_i$  are diagonal matrices.

# Глубокие рекуррентные сети (DLSTM)

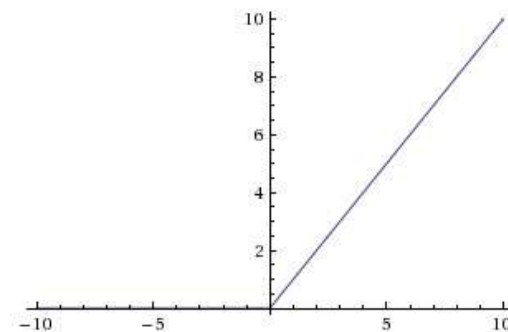


Amodei (2015) *Deep Speech 2: End-to-End Speech Recognition in English and Mandarin*

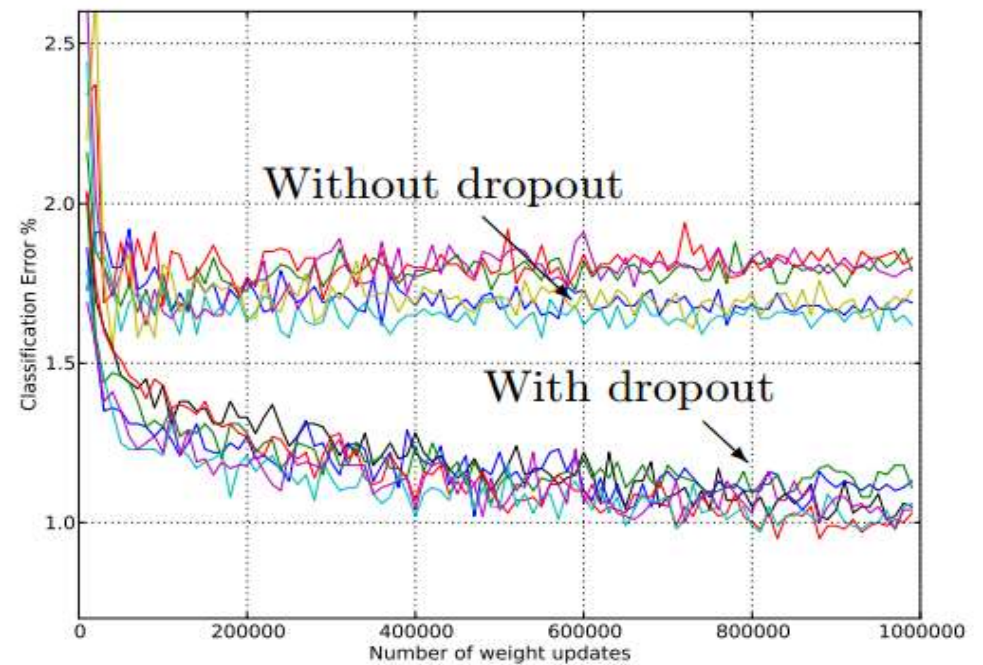
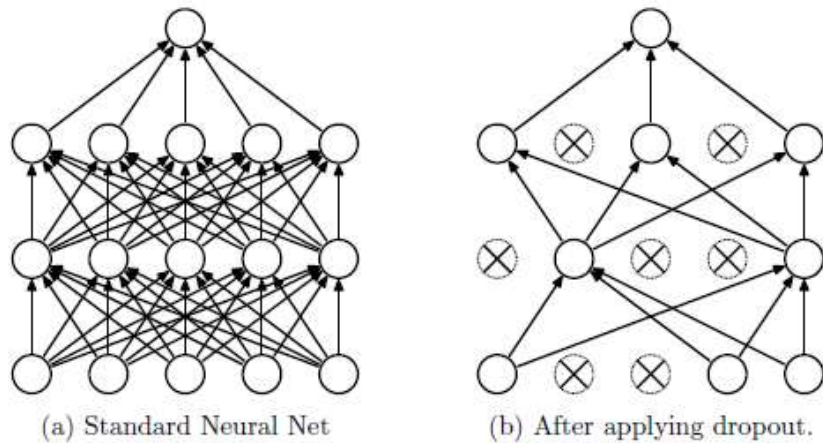


# Регуляризация обучения

- **ReLU** (Nair, [2010](#))
- **Dropout** (Hinton, [2012](#))
- **Batch normalization** (Ioffe, [2015](#))



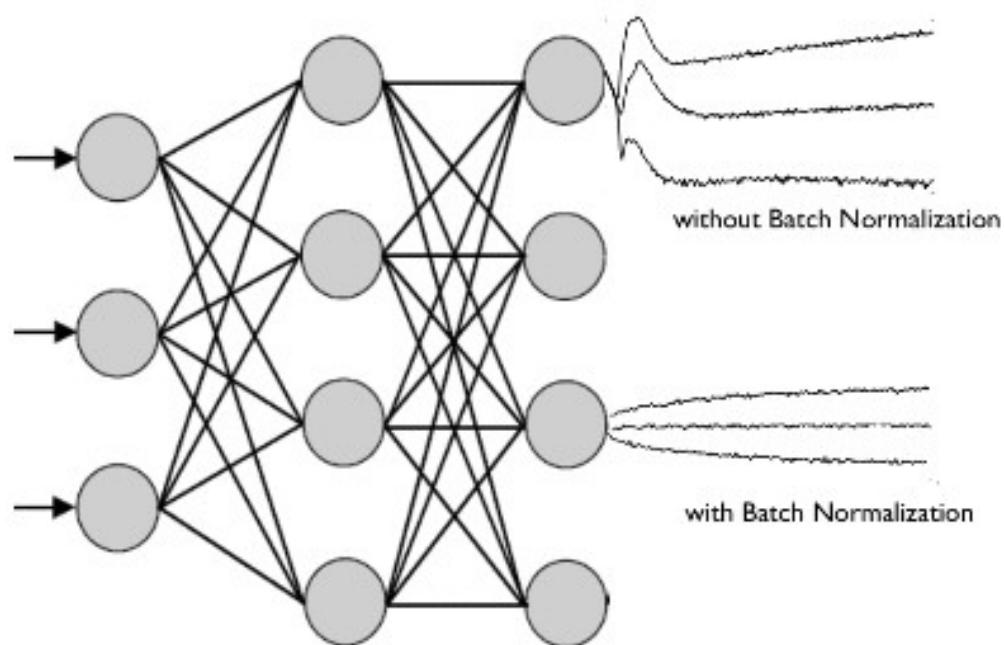
# Dropout



Hinton ([2012](#)) *Improving neural networks by preventing co-adaptation of feature detectors*

# Batch normalization

Ioffe (2015) *Batch normalization: Accelerating deep network training by reducing internal covariate shift*

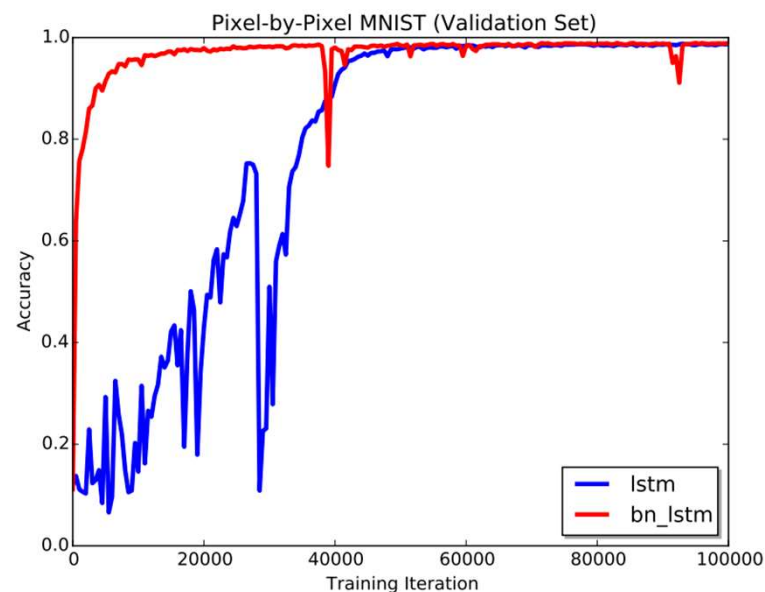


$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i$$

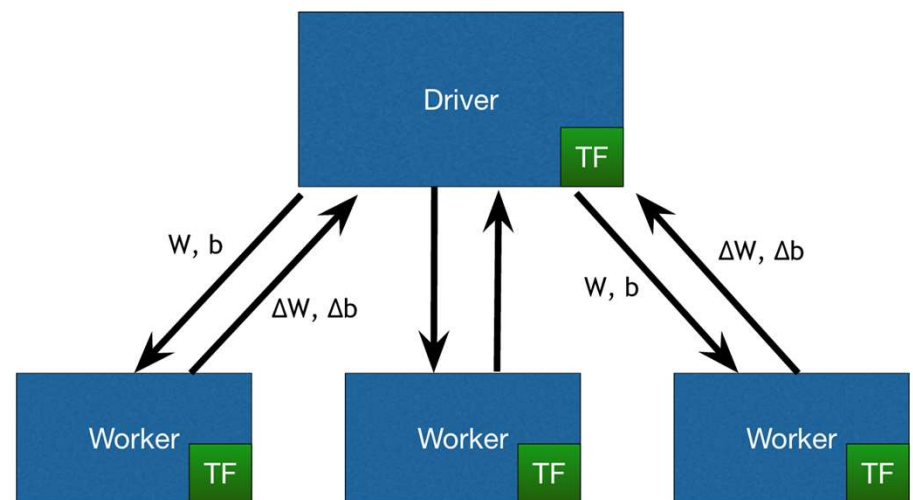
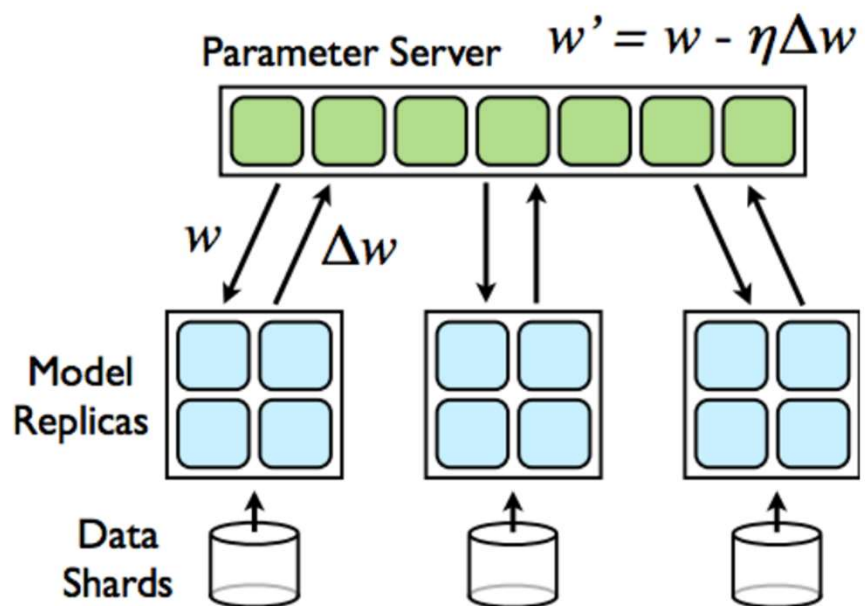
$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i)$$



# Parallelizing and distributing SGD Hogwild!



# Downpour SGD

Downpour SGD is an asynchronous variant of SGD that was used by Dean et al. [\[4\]](#) in their DistBelief framework (predecessor to TensorFlow) at Google. It runs multiple replicas of a model in parallel on subsets of the training data. These models send their updates to a parameter server, which is split across many machines. Each machine is responsible for storing and updating a fraction of the model's parameters. However, as replicas don't communicate with each other e.g. by sharing weights or updates, their parameters are continuously at risk of diverging, hindering convergence.