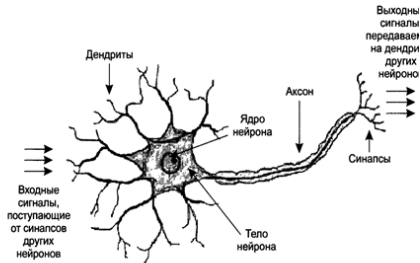


# Neural Nets (NN), с элементами Deep Learning

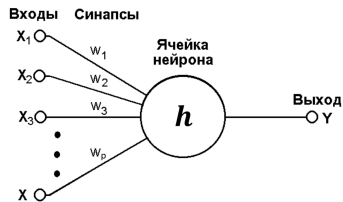
Е. Ларин, Ф. Ежов, И. Кононыхин

Санкт-Петербургский государственный университет  
Прикладная математика и информатика  
Вычислительная стохастика и статистические модели

# Математическая модель нейрона



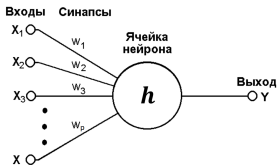
## Формальный нейрон



$$h = \sum_i x_i * w_i \quad y = f(h)$$

# Перцептрон

## Формальный нейрон



$$h = \sum_i x_i * w_i \quad y = f(h)$$

Перцептрон(perceptron)  
— одна из простейших архитектур ANN, придуманная Фрэнком Розенблаттом в 1957 году.

- $p$  — количество признаков
- $\{x_1, \dots, x_p\}$  — индивид.
- $\{w_1, \dots, w_p\}$  — веса нейрона.

- 

$$f(h) = \begin{cases} 0 & h < 0 \\ 1 & h \geq 0 \end{cases}$$

— функция активации.

- $y$  — выход сети.

# Перцептрон: Проблемы

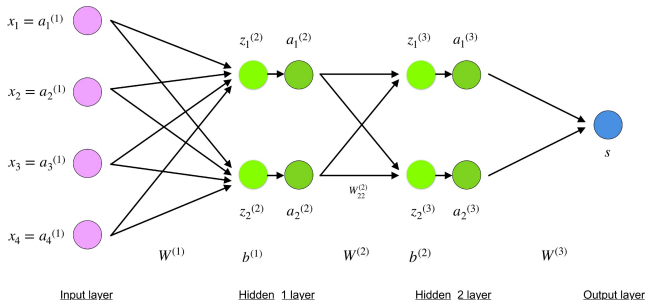
Перцептрон может решать только линейно разделимые задачи.

## Проблема

Перцептрон не способен решить некоторые тривиальные задачи, например задачу классификации на основе исключающего "ИЛИ" (Exclusive OR - XOR).

# MultiLayer Perceptron (MLP)

Каждый слой кроме выходного включает нейрон смещения и полностью связан со следующим слоем.



# MultiLayer Perceptron (MLP): Теорема Цыбенко

Функция  $\sigma(z)$  — сигмоида, если  $\lim_{z \rightarrow -\infty} \sigma(z) = 0$  и  $\lim_{z \rightarrow +\infty} \sigma(z) = 1$ .

## Теорема Цыбенко

Если функция активации  $\sigma(z)$  — непрерывная сигмоида, то для любой непрерывной на  $[0, 1]^p$  функции  $f(x)$  существуют такие значения параметров  $w_h \in \mathbb{R}^p$ ,  $w_0 \in \mathbb{R}$ ,  $\alpha_h \in \mathbb{R}$ , что двухслойная сеть

$$a(x) = \sum_{h=1}^H \alpha_h \sigma(\langle x, w_h \rangle + w_0)$$

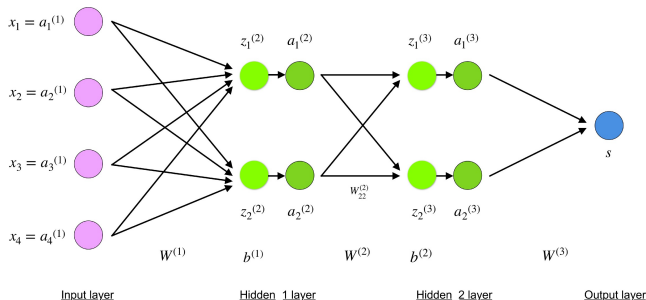
равномерно приближает  $f(x)$  с любой точностью  $\varepsilon$ :

$$|a(x) - f(x)| < \varepsilon, \text{ для всех } x \in [0, 1]^p$$

# MultiLayer Perceptron (MLP): Функции представимые нейросетью

- Двухслойная сеть в  $\mathbb{R}^p$  позволяет отделить произвольный выпуклый многогранник.
- Трехслойная сеть в  $\mathbb{R}^p$  позволяет отделить произвольную многогранную область, не обязательно выпуклую, и даже не обязательно связную.
- С помощью линейных операций и одной нелинейной функции активации можно приблизить любую непрерывную функцию с любой желаемой точностью.

# Обучение MLP



Минимизация средних потерь на обучающей выборке:

$$Q(w) := \frac{1}{n} \sum_{i=1}^n \mathcal{C}_i(w) \rightarrow \min_w$$

$\mathcal{C}$  — любая дифференцируемая функция потерь.



# Алгоритм SG

**Задача минимизации:**  $Q(w) := \frac{1}{n} \sum_{i=1}^n \mathcal{C}_i(w) \rightarrow \min_w$

**Вход:** выборка  $X^n$ ; «скорость» обучения  $\epsilon$ ; параметр  $\lambda$ ;

**Выход:** оптимизированные значение весов  $w$ .

- ❶ инициализация весов  $w$  и текущую оценку  $Q(w)$
- ❷ **повторять**
- ❸ выбрать объект  $x_i$  из  $X^n$ ;
- ❹ вычислить потерю  $\mathcal{C}_i := \mathcal{C}_i(w)$
- ❺ градиентный шаг:  $w := w - \epsilon \mathcal{C}'_i(w)$
- ❻ оценить значение функционала:  $Q := (1 - \lambda)Q + \lambda \mathcal{C}_i$
- ❼ **пока** значение  $Q$  и/или весов  $w$  не стабилизируются

# Обучение MLP: forward-propagation

$$x = a^{(1)} \quad \text{Input layer}$$

$$z^{(2)} = W^{(1)}x + b^{(1)} \quad \text{neuron value at Hidden}_1 \text{ layer}$$

$$a^{(2)} = f(z^{(2)}) \quad \text{activation value at Hidden}_1 \text{ layer}$$

$$z^{(3)} = W^{(2)}a^{(2)} + b^{(2)} \quad \text{neuron value at Hidden}_2 \text{ layer}$$

$$a^{(3)} = f(z^{(3)}) \quad \text{activation value at Hidden}_2 \text{ layer}$$

$$s = W^{(3)}a^{(3)} \quad \text{Output layer}$$

$$C = \text{cost}(s, y)$$

$s$  — полученная оценка с помощью нейронной сети.  $y$  — правильный ответ.  $C$  — значение ошибки между  $s$  и  $y$  по какой-нибудь метрики (например MSE).

# Обучение MLP: back-propagation

$$\frac{\partial C}{\partial w_{jk}^l} = \frac{\partial C}{\partial z_j^l} \frac{\partial z_j^l}{\partial w_{jk}^l} \quad \text{chain rule}$$

$$z_j^l = \sum_{k=1}^m w_{jk}^l a_k^{l-1} + b_j^l \quad \text{by definition}$$

$m$  – number of neurons in  $l-1$  layer

$$\frac{\partial z_j^l}{\partial w_{jk}^l} = a_k^{l-1} \quad \text{by differentiation (calculating derivative)}$$

$$\frac{\partial C}{\partial w_{jk}^l} = \frac{\partial C}{\partial z_j^l} a_k^{l-1} \quad \text{final value}$$

## Обучение MLP: процесс изменения весов

*while (termination condition not met)*

$$w := w - \epsilon \frac{\partial C}{\partial w}$$

$$b := b - \epsilon \frac{\partial C}{\partial b}$$

*end*

# back-propagation algorithm: в итоге

## Преимущества:

- быстрое вычисление градиента
- обобщение на любые функции активации, функции потерь и количество слоев
- возможность потокового обучения
- возможность распараллеливания

## Недостатки:

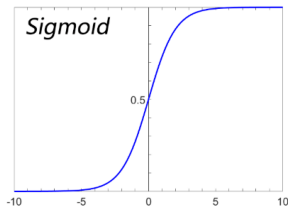
- медленная сходимость
- застревания в локальных экстремумах
- проблема переобучения

# Методы решения задачи минимизации

- Stochastic Gradient (SG)
- Метод накопления импульса (Momentum)
- NAG (Nesterov's accelerated gradient)
- RMSProp (running mean square)
- AdaDelta (adaptive learning rate)
- Adam (adaptive momentum)
- Nadam (Nesterov-accelerated adaptive momentume)

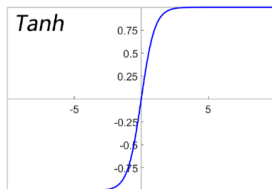
# Функции активации

Функции активации добавляют в модель нелинейность.



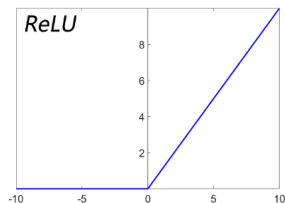
$$f(x) = \frac{1}{1 + e^{-x}}$$

(a)



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

(b)



$$f(x) = \max(0, x)$$

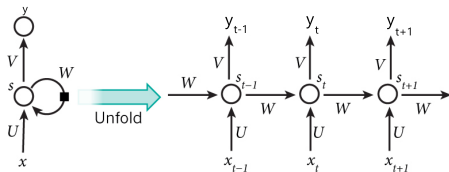
(c)

# Recurrent neural network(RNN)

$x_t$  — входной вектор в момент времени  $t$

$s_t$  — вектор скрытого слоя в момент времени  $t$

$y_t$  — выходной вектор (в некоторых случаях  $y_t \equiv s_t$ )



Задача:

$$\sum_{t=0}^T \mathcal{C}_t(U, V, W) \rightarrow \min_{U, V, W}$$

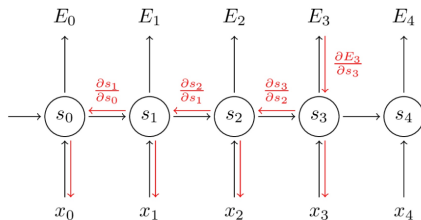
$\mathcal{C}_t(U, V, W) = \mathcal{C}(y_t(U, V, W))$  — потеря от предсказания  $y_t$ .



# Обучение RNN

## Back-propagation Through Time (BPTT)

$$\frac{\partial \mathcal{C}_t}{\partial W} = \frac{\partial \mathcal{C}_t}{\partial y_t} \frac{\partial y_t}{\partial s_t} \sum_{k=0}^t \left( \prod_{i=k+1}^t \frac{\partial s_i}{\partial s_{i-1}} \right) \frac{\partial h_k}{\partial W}$$



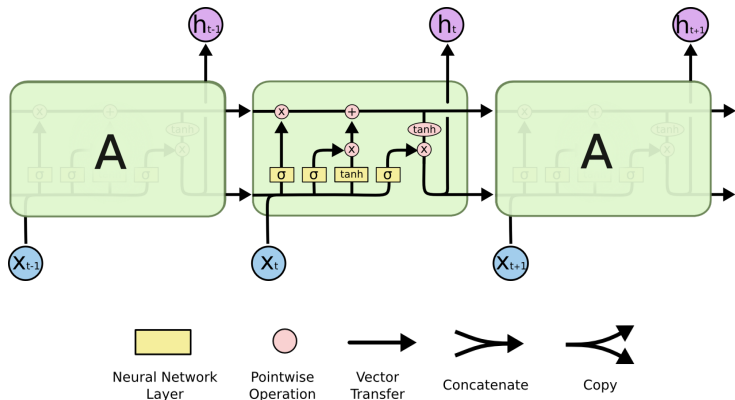
# Где используется RNN?

- Прогнозирование временных рядов
- Управление технологическими процессами
- Классификация текстов или их фрагментов
- Анализ тональности документа / предложений / слов
- Машинный перевод
- Распознавание речи
- Синтез речи
- Синтез ответов на вопросы, разговорный интеллект
- Генерация подписей к изображениям
- Генерация рукописного текста
- Интерпретация генома и другие задачи биоинформатики

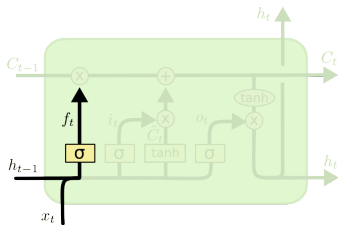
# Long short-term memory (LSTM)

**Мотивация LSTM:** сеть должна долго помнить контекст, какой именно — сеть должна выучить сама.


Вводится  $C_t$  — вектор состояния сети в момент  $t$ .



# Объясняя LSTM (часть 1)



$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

  
Neural Network  
Layer

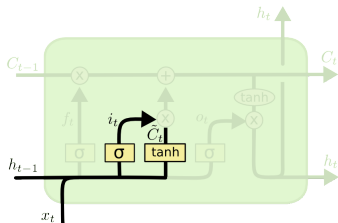
  
Pointwise  
Operation

  
Vector  
Transfer

  
Concatenate

  
Copy

## Объясняя LSTM (часть 2)



$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

Neural Network  
Layer

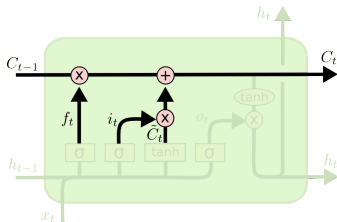
Pointwise  
Operation

Vector  
Transfer


Concatenate

Copy

# Объясняя LSTM (часть 3)



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

  
Neural Network  
Layer

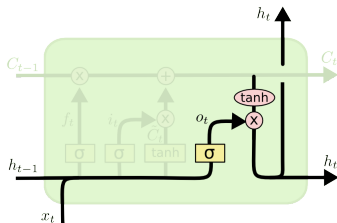
  
Pointwise  
Operation

  
Vector  
Transfer

  
Concatenate

  
Copy

# Объясняя LSTM (часть 4)



$$o_t = \sigma(W_o [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh(C_t)$$

Neural Network  
Layer

Pointwise  
Operation

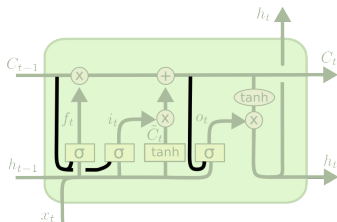
Vector  
Transfer

Concatenate

Copy

# Модификации LSTM

## LSTM с «замочными скважинами» (peepholes)



$$f_t = \sigma(W_f \cdot [C_{t-1}, h_{t-1}, x_t] + b_f)$$

$$i_t = \sigma(W_i \cdot [C_{t-1}, h_{t-1}, x_t] + b_i)$$

$$o_t = \sigma(W_o \cdot [C_t, h_{t-1}, x_t] + b_o)$$

Все фильтры «подглядывают» вектор состояния  $C_{t-1}$  или  $C_t$ .

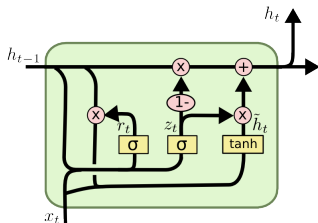
Увеличивается число параметров модели.

Замочную скважину можно использовать не для всех фильтров.



# Модификации LSTM

## LSTM Gated recurrent units (GRU)



$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

Используется только состояние  $h_t$ , вектор  $C_t$  не вводится.

Фильтр обновления (update gate) вместо входного и забывающего.

Фильтр перезагрузки (reset gate) решает, какую часть памяти нужно перенести дальше с прошлого шага.

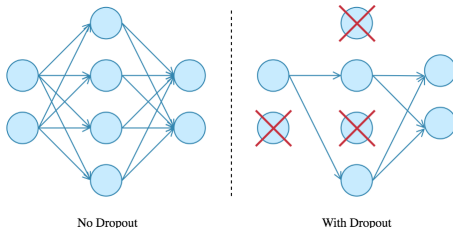
# Dropout

**Обучение:** обнуляем выход  $h$ -ого нейрона  $\ell$ -го слоя с вероятностью  $p_\ell$ .

$$a_{ih}^{\ell+1} = \xi_h^\ell f_h^\ell \left( \sum_j w_{jh} a_{ij}^\ell \right), P(\xi_h^\ell = 0) = p_\ell$$

**Применение:** вводим поправку

$$a_{ih}^{\ell+1} = (1 - p_\ell) f_h^\ell \left( \sum_j w_{jh} a_{ij}^\ell \right)$$



# Dropout

## Inverted Dropout:

$$a_{ih}^{\ell+1} = \frac{1}{(1 - p_\ell)} \xi_h^\ell f_h^\ell \left( \sum_j w_{jh} a_{ij}^\ell \right)$$

- регуляризация: из всех сетей выбираем более устойчивую к утрате  $pN$  нейронов.
- сокращаем переобучение, заставляя разные части сети решать одну и ту же исходную задачу.

# Batch Normalization

$B = x_i$  — пакеты (mini-batch) данных. Усреднение градиентов  $\mathcal{C}(w)$  по пакету ускоряет сходимость.

$B^\ell = \{u_i^\ell\}$  — векторы объектов  $x_i$  на выходе  $\ell$ -го слоя.

## Batch Normalization:

- 1 Нормировать каждую  $j$ -ю компоненту вектора  $u_i^\ell$  по пакету:

$$\hat{u}_{ij}^\ell = \frac{u_{ij}^\ell - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}; \quad \nu_j = \frac{1}{|B|} \sum_{x_i \in B} u_{ij}^\ell; \quad \sigma_j^2 = \frac{1}{|B|} \sum_{x_i \in B} (u_{ij}^\ell - \nu_j)^2.$$

- 2 Добавить линейный слой с настраиваемыми весами:

$$\tilde{u}_{ij}^\ell = \gamma_j^\ell \hat{u}_{ij}^\ell + \beta_j^\ell$$

- 3 Параметры  $\gamma_j^\ell$  и  $\beta_j^\ell$  настраиваются с помощью back-propagation.