

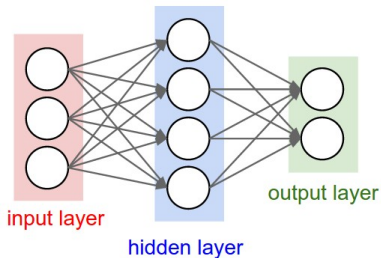
# Лекция 2. Обучение нейронных сетей. Backpropagation

Глубинное обучение

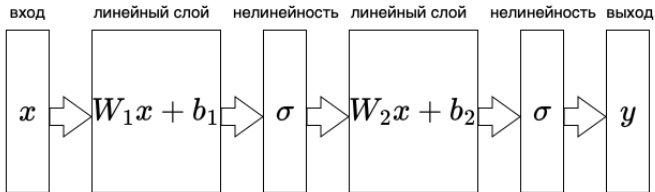
---

Антон Кленицкий

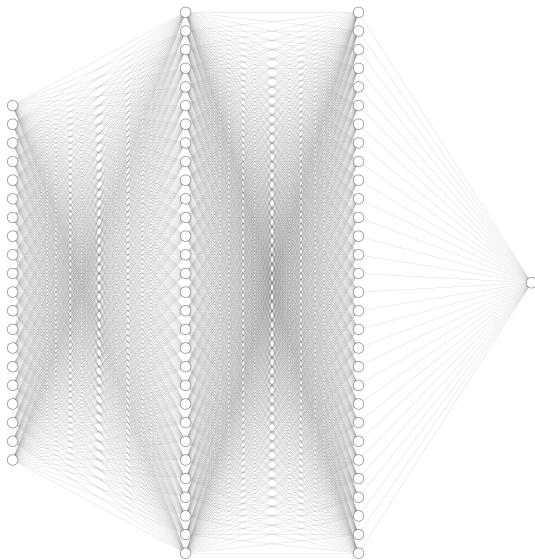
# Multilayer Perceptron



$$y = \sigma(W_2\sigma(W_1x+b_1)+b_2)$$



В реальной жизни скорее так



## Теорема об универсальной аппроксимации:

Любую непрерывную функцию можно с любой точностью приблизить нейросетью с одним скрытым слоем с сигмоидной функцией активации.

- Hornik "Multilayer feedforward networks are universal approximators"(1989)
- Cybenko "Approximation by superpositions of a sigmoidal function"(1989)

Но!

- Может понадобится очень много (экспоненциально много) нейронов
- Неизвестно, сможем ли мы обучить нашими методами такую нейросеть
- Неизвестно, как это будет обобщаться на новые данные

# Почему же "глубокое" обучение?

## Глубокие нейронные сети

- $> 1$  скрытого слоя
- Строят иерархию признаков
- Имеют большую выразительность при том же числе нейронов
- Работают на практике



Image credit

## Convolutional neural networks (CNN) - Сверточные сети

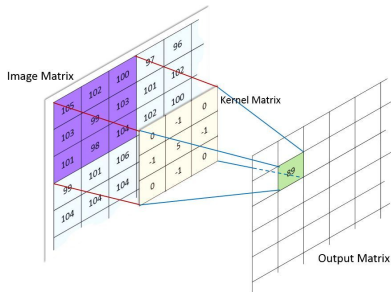


Image credit

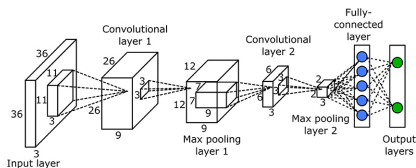


Image credit

- Могут быть действительно глубокими (100+ слоев)
- Изображения, звук, последовательности

# Основные архитектуры нейросетей

Recurrent neural networks (RNN) - Рекуррентные сети

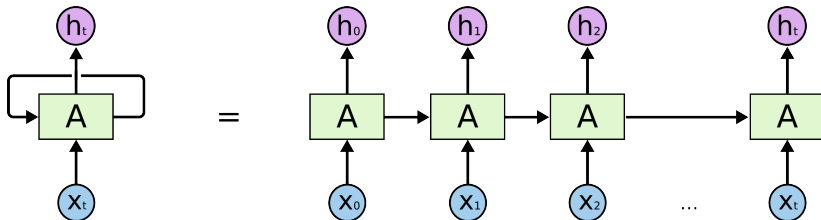


Image credit

Для работы с последовательностями



# Основные архитектуры нейросетей

## Transformers

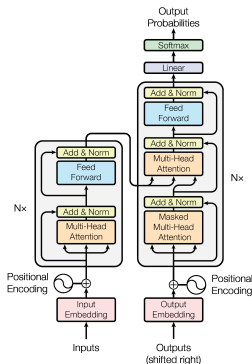


Image credit

NLP, последовательности, изображения, таблицы

## Autoencoders

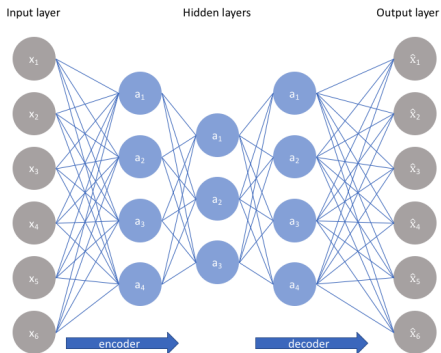
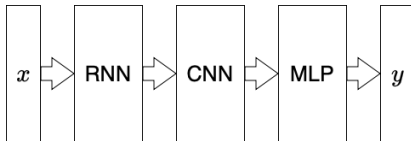
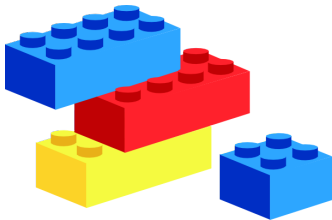


Image credit

# Модель можно составлять из разных блоков



Как обучать нейросети?

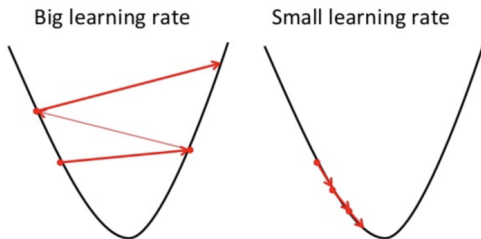
---

# Gradient descent

$$L(\theta) = \frac{1}{N} \sum_{i=1}^N L(y_i, f(x_i, \theta)) \rightarrow \min_{\theta}$$

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \nabla_{\theta} L(\theta^{(t)})$$

$$\nabla_{\theta} L(\theta) = \left( \frac{\partial L}{\partial \theta_1}, \frac{\partial L}{\partial \theta_2}, \dots \right)$$

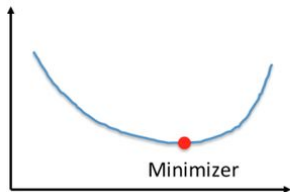


# Convex vs non-convex optimization

Выпуклая оптимизация

Невыпуклая оптимизация

**Convex**



**Non-Convex**

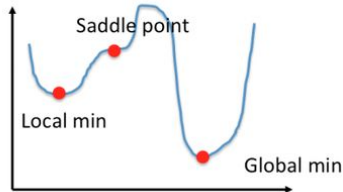
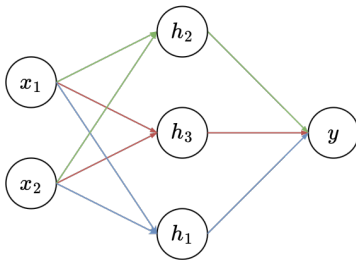
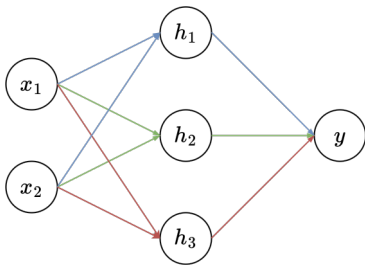


Image credit

Нейросети - сложные функции с огромным количеством локальных минимумов

# Локальные минимумы у нейросетей

Изменится ли выход нейросети при перемешивании нейронов в скрытом слое?



(Batch) gradient descent

$$\theta_{t+1} = \theta_t - \alpha \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} L(y_i, f(x_i, \theta))$$



# Stochastic Gradient Descent

Stochastic gradient descent (SGD) - обновляем веса после каждого примера

$$\theta_{t+1} = \theta_t - \alpha \nabla_{\theta} L(y_i, f(x_i, \theta))$$

- Несмещенная оценка полного градиента
- Если берем примеры из обучающей выборки в случайном порядке!

Mini-batch stochastic gradient descent - обновляем веса после батча из  $B$  примеров

$$\theta_{t+1} = \theta_t - \alpha \frac{1}{B} \sum_{i=1}^B \nabla_{\theta} L(y_i, f(x_i, \theta))$$

- Можно эффективно использовать матричные вычисления

# Stochastic Gradient Descent

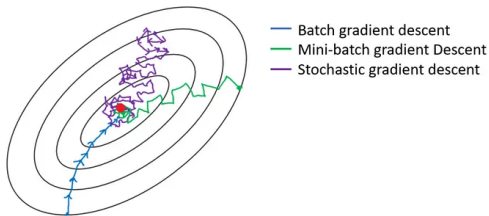


Image credit

- Шумная оценка полного градиента
- Может проскочить мимо неудачного локального минимума или седловой точки

# Wide vs sharp minimum

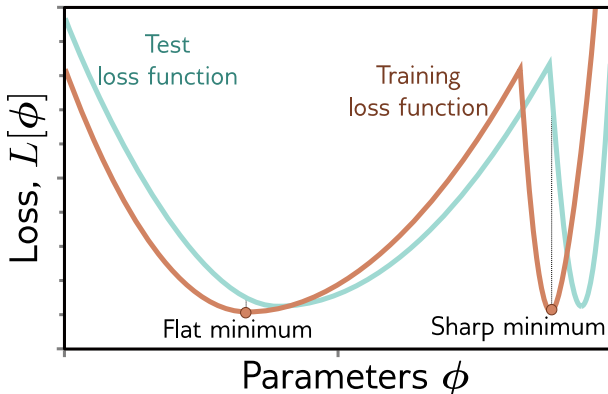
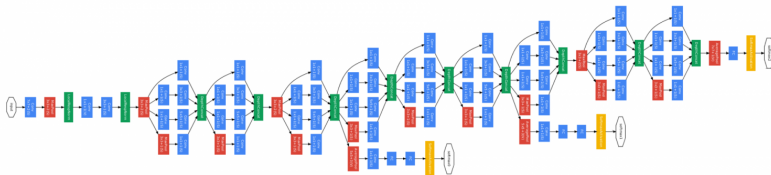


Image credit

# Как считать градиенты?

Аналитически?



Но модель может быть очень большая и сложная :)

# Как считать градиенты?

Численно?

Конечные разности

$$\frac{\partial f(x, \theta)}{\partial \theta} \approx \frac{f(x, \theta + \epsilon) - f(x, \theta - \epsilon)}{2\epsilon}$$

- Приближенный метод
- Очень долго - для каждого параметра по отдельности
- Можно использовать для проверки

# Backpropagation

---

Алгоритм обратного распространения ошибки

## Yes you should understand backprop



Andrej Karpathy · Follow

7 min read · Dec 19, 2016

[https://karpathy.medium.com/  
yes-you-should-understand-backprop-e2f06eab496b](https://karpathy.medium.com/yes-you-should-understand-backprop-e2f06eab496b)

# Правило дифференцирования сложной функции

$$f = f(g(x))$$

$$\frac{df}{dx} = \frac{df}{dg} \frac{dg}{dx}$$

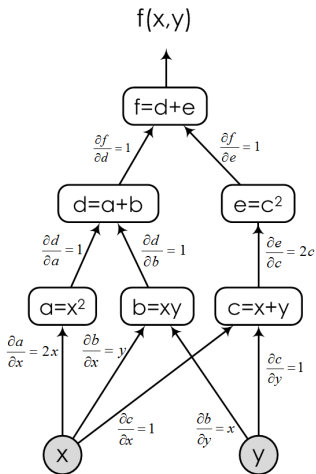
$$f = f(g_1(x), g_2(x), \dots, g_k(x))$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial g_1} \frac{\partial g_1}{\partial x} + \dots + \frac{\partial f}{\partial g_k} \frac{\partial g_k}{\partial x} = \sum_{i=1}^k \frac{\partial f}{\partial g_i} \frac{\partial g_i}{\partial x}$$



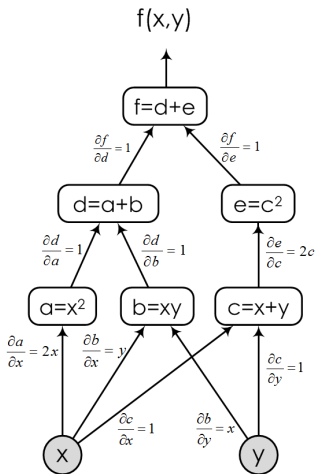
# Простой пример

$$f(x) = x^2 + xy + (x + y)^2$$



Граф вычислений

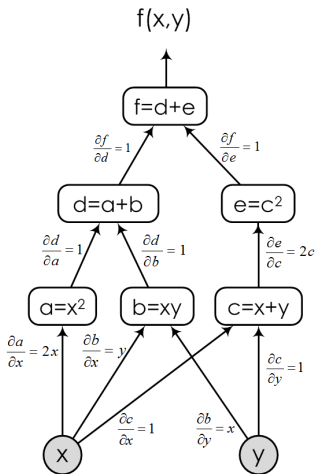
# Простой пример



$$f(x) = x^2 + xy + (x + y)^2$$

$$\frac{\partial f}{\partial d} = 1; \frac{\partial f}{\partial e} = 1$$

# Простой пример



$$f(x) = x^2 + xy + (x + y)^2$$

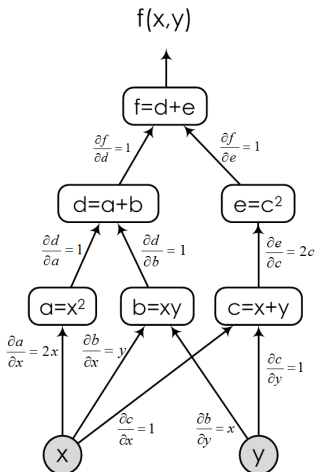
$$\frac{\partial f}{\partial d} = 1; \frac{\partial f}{\partial e} = 1$$

$$\frac{\partial f}{\partial a} = \frac{\partial f}{\partial d} \frac{\partial d}{\partial a} = 1$$

$$\frac{\partial f}{\partial b} = \frac{\partial f}{\partial d} \frac{\partial d}{\partial b} = 1$$

$$\frac{\partial f}{\partial c} = \frac{\partial f}{\partial e} \frac{\partial e}{\partial c} = 2c$$

# Простой пример



Граф вычислений

$$f(x) = x^2 + xy + (x + y)^2$$

$$\frac{\partial f}{\partial d} = 1; \frac{\partial f}{\partial e} = 1$$

$$\frac{\partial f}{\partial a} = \frac{\partial f}{\partial d} \frac{\partial d}{\partial a} = 1$$

$$\frac{\partial f}{\partial b} = \frac{\partial f}{\partial d} \frac{\partial d}{\partial b} = 1$$

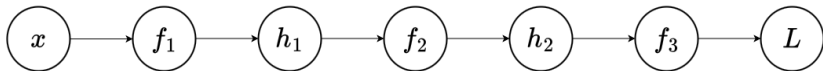
$$\frac{\partial f}{\partial c} = \frac{\partial f}{\partial e} \frac{\partial e}{\partial c} = 2c$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial a} \frac{\partial a}{\partial x} + \frac{\partial f}{\partial b} \frac{\partial b}{\partial x} + \frac{\partial f}{\partial c} \frac{\partial c}{\partial x} = 2x + y + 2(x + y)$$

$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial b} \frac{\partial b}{\partial y} + \frac{\partial f}{\partial c} \frac{\partial c}{\partial y} = x + 2(x + y)$$

$$\hat{y} = w_3 \sigma(w_2 \sigma(w_1 x + b_1) + b_2) + b_3; \quad L(y, \hat{y}) = (y - \hat{y})^2$$

$$\hat{y} = w_3\sigma(w_2\sigma(w_1x + b_1) + b_2) + b_3; \quad L(y, \hat{y}) = (y - \hat{y})^2$$



$$f_1 = w_1x + b_1$$

$$h_1 = \sigma(f_1)$$

$$f_2 = w_2h_1 + b_2$$

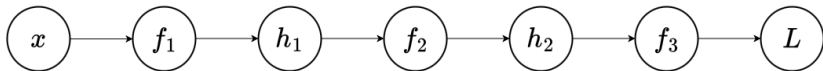
$$h_2 = \sigma(f_2)$$

$$f_3 = w_3h_2 + b_3$$

$$L = (f_3 - y)^2$$

# Простая нейросеть

$$\hat{y} = w_3\sigma(w_2\sigma(w_1x + b_1) + b_2) + b_3; \quad L(y, \hat{y}) = (y - \hat{y})^2$$



$$\frac{\partial L}{\partial f_3} = 2(f_3 - y)$$

$$f_1 = w_1x + b_1$$

$$h_1 = \sigma(f_1)$$

$$f_2 = w_2h_1 + b_2$$

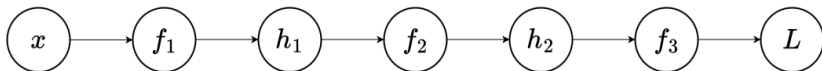
$$h_2 = \sigma(f_2)$$

$$f_3 = w_3h_2 + b_3$$

$$L = (f_3 - y)^2$$

# Простая нейросеть

$$\hat{y} = w_3\sigma(w_2\sigma(w_1x + b_1) + b_2) + b_3; \quad L(y, \hat{y}) = (y - \hat{y})^2$$



$$\frac{\partial L}{\partial f_3} = 2(f_3 - y)$$

$$f_1 = w_1x + b_1$$

$$h_1 = \sigma(f_1)$$

$$f_2 = w_2h_1 + b_2$$

$$h_2 = \sigma(f_2)$$

$$f_3 = w_3h_2 + b_3$$

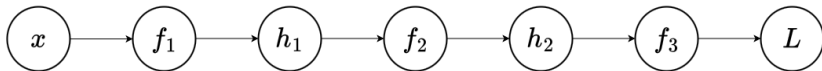
$$L = (f_3 - y)^2$$

$$\frac{\partial L}{\partial h_2} = \left( \frac{\partial L}{\partial f_3} \right) \frac{\partial f_3}{\partial h_2}$$



# Простая нейросеть

$$\hat{y} = w_3\sigma(w_2\sigma(w_1x + b_1) + b_2) + b_3; \quad L(y, \hat{y}) = (y - \hat{y})^2$$



$$\frac{\partial L}{\partial f_3} = 2(f_3 - y)$$

$$f_1 = w_1x + b_1$$

$$h_1 = \sigma(f_1)$$

$$f_2 = w_2h_1 + b_2$$

$$h_2 = \sigma(f_2)$$

$$f_3 = w_3h_2 + b_3$$

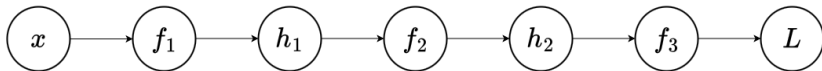
$$L = (f_3 - y)^2$$

$$\frac{\partial L}{\partial h_2} = \left( \frac{\partial L}{\partial f_3} \right) \frac{\partial f_3}{\partial h_2}$$

$$\frac{\partial L}{\partial f_2} = \left( \frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial h_2} \right) \frac{\partial h_2}{\partial f_2}$$

# Простая нейросеть

$$\hat{y} = w_3 \sigma(w_2 \sigma(w_1 x + b_1) + b_2) + b_3; \quad L(y, \hat{y}) = (y - \hat{y})^2$$



$$\frac{\partial L}{\partial f_3} = 2(f_3 - y)$$

$$f_1 = w_1 x + b_1$$

$$h_1 = \sigma(f_1)$$

$$f_2 = w_2 h_1 + b_2$$

$$h_2 = \sigma(f_2)$$

$$f_3 = w_3 h_2 + b_3$$

$$L = (f_3 - y)^2$$

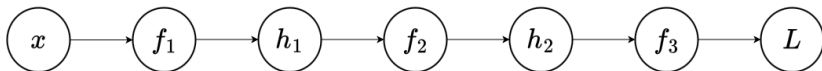
$$\frac{\partial L}{\partial h_2} = \left( \frac{\partial L}{\partial f_3} \right) \frac{\partial f_3}{\partial h_2}$$

$$\frac{\partial L}{\partial f_2} = \left( \frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial h_2} \right) \frac{\partial h_2}{\partial f_2}$$

$$\frac{\partial L}{\partial h_1} = \left( \frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial h_2} \frac{\partial h_2}{\partial f_2} \right) \frac{\partial f_2}{\partial h_1}$$

# Простая нейросеть

$$\hat{y} = w_3\sigma(w_2\sigma(w_1x + b_1) + b_2) + b_3; \quad L(y, \hat{y}) = (y - \hat{y})^2$$



$$\frac{\partial L}{\partial f_3} = 2(f_3 - y)$$

$$f_1 = w_1x + b_1$$

$$h_1 = \sigma(f_1)$$

$$f_2 = w_2h_1 + b_2$$

$$h_2 = \sigma(f_2)$$

$$f_3 = w_3h_2 + b_3$$

$$L = (f_3 - y)^2$$

$$\frac{\partial L}{\partial h_2} = \left( \frac{\partial L}{\partial f_3} \right) \frac{\partial f_3}{\partial h_2}$$

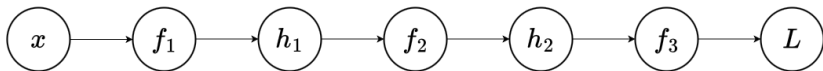
$$\frac{\partial L}{\partial f_2} = \left( \frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial h_2} \right) \frac{\partial h_2}{\partial f_2}$$

$$\frac{\partial L}{\partial h_1} = \left( \frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial h_2} \frac{\partial h_2}{\partial f_2} \right) \frac{\partial f_2}{\partial h_1}$$

$$\frac{\partial L}{\partial f_1} = \left( \frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial h_2} \frac{\partial h_2}{\partial f_2} \frac{\partial f_2}{\partial h_1} \right) \frac{\partial h_1}{\partial f_1}$$

# Простая нейросеть

$$\hat{y} = w_3\sigma(w_2\sigma(w_1x + b_1) + b_2) + b_3; \quad L(y, \hat{y}) = (y - \hat{y})^2$$



$$f_1 = w_1x + b_1$$

$$h_1 = \sigma(f_1)$$

$$f_2 = w_2h_1 + b_2$$

$$h_2 = \sigma(f_2)$$

$$f_3 = w_3h_2 + b_3$$

$$L = (f_3 - y)^2$$

$$\frac{\partial L}{\partial w_k} = \left( \frac{\partial L}{\partial f_k} \right) \frac{\partial f_k}{\partial w_k} = \frac{\partial L}{\partial f_k} \cdot h_{k-1}$$

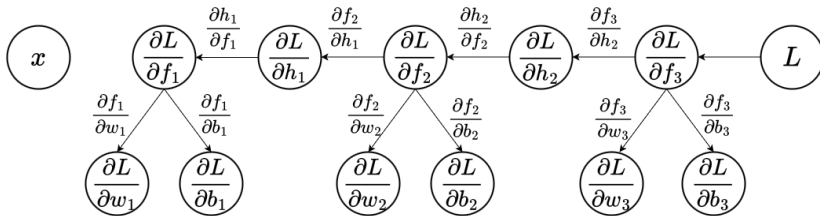
$$\frac{\partial L}{\partial b_k} = \left( \frac{\partial L}{\partial f_k} \right) \frac{\partial f_k}{\partial b_k} = \frac{\partial L}{\partial f_k} \cdot 1$$

# Простая нейросеть

Forward pass



Backward pass



# Backpropagation summary

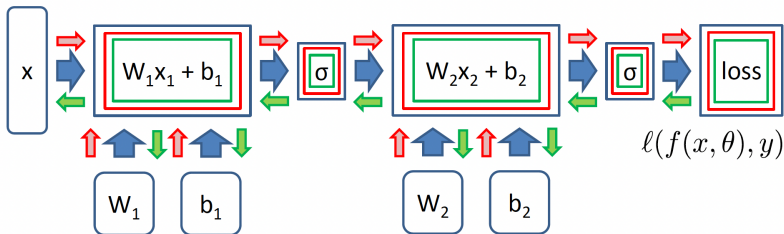


Image credit

- Нейросеть - граф вычислений
- Градиенты распространяются в обратную сторону по графу
- Умножаем на производную выхода по входу для каждого узла графа

В процессе обратного распространения по сети градиенты могут "затухать" или "взрываться"

$$\frac{\partial L}{\partial f_1} = \left( \frac{\partial L}{\partial f_3} \frac{\partial f_3}{\partial h_2} \frac{\partial h_2}{\partial f_2} \frac{\partial f_2}{\partial h_1} \right) \frac{\partial h_1}{\partial f_1}$$

# Let the gradient flow

Главное, чтобы градиент хорошо "протекал" по нейросети

Многое было придумано ради этого

- Хорошие функции активации
- Хорошая инициализация весов
- Skip residual connections
- Специальные архитектуры (LSTM/GRU)
- ...



# DL Frameworks

---

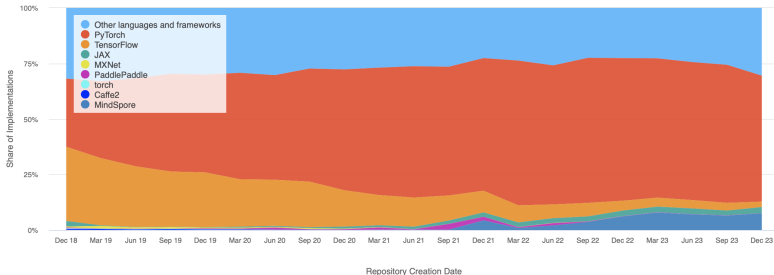
Автоматическое дифференцирование - расчет производных по графу вычислений



Image credit

<https://paperswithcode.com/trends>

Paper Implementations grouped by framework



## В следующий раз

- Функции активации
- Функции потерь
- Инициализация весов