

Лекция 1. Введение

Денис Деркач, Дмитрий Тарасов



LAMBDA • HSE

19 января 2026 года

Параметры курса



Лекции



Денис Деркач

- Магистратура (СПбГУ), 2007, PhD (Париж 11), 2010.
- Заведующий лабораторией методов анализа больших данных (LAMBDA)

Интересы

- Генеративное моделирование физических процессов с помощью машинного обучения.
- Оптимизация и настройка инженерного оборудования с помощью ML.
- Объединение DL и классических моделей.

Идея курса

DL – довольно размытый термин, потому зафиксируем некоторые детали.

К концу курса вы будете:

- ▶ Понимать базовые принципы работы глубокого обучения
- ▶ Использовать на практике полученные знания.
- ▶ Выбирать согласно данным основные архитектуры нейронных сетей - свёрточные, рекуррентные, трансформеры
- ▶ Использовать в основных приложения - Computer Vision, NLP
- ▶ Применять базовые архитектуры генеративных моделей в стандартных постановках задачи.

Организационные вопросы

- ▶ Материалы курса: <http://github.com/fintech-dl-hse/course>
- ▶ Формула оценки: $0.8 \cdot \sum_i O_{hw_i} + 0.1 \cdot O_{exam3} + 0.1 \cdot O_{exam4}$
- ▶ После занятия в боте будет открываться квиз с открытыми вопросами по темам занятия. Вопросы несложные. Прохождение квизов будет допуском к экзаменам отдельно и за 3 и за 4 модуль.
- ▶ Telegram чат
- ▶ Форма обратной связи после занятий
- ▶ Задавайте вопросы!



Литература

- ▶ Глубокое обучение. Погружение в мир нейронных сетей. С. Николенко
- ▶ А. Кадулин, Е. Архангельская. • Глубокое обучение.
- ▶ Я. Гудфеллоу, Й. Бенджио, А. Курвилль. Английская версия:
<http://www.deeplearningbook.org/>.
- ▶ Understanding deep learning. S. Prince <http://udlbook.github.io/udlbook/>
- ▶ Dive into Deep Learning. A. Zhang, Z. Lipton, M. Li, A. Smola <http://d2l.ai/>

Дополнительные источники

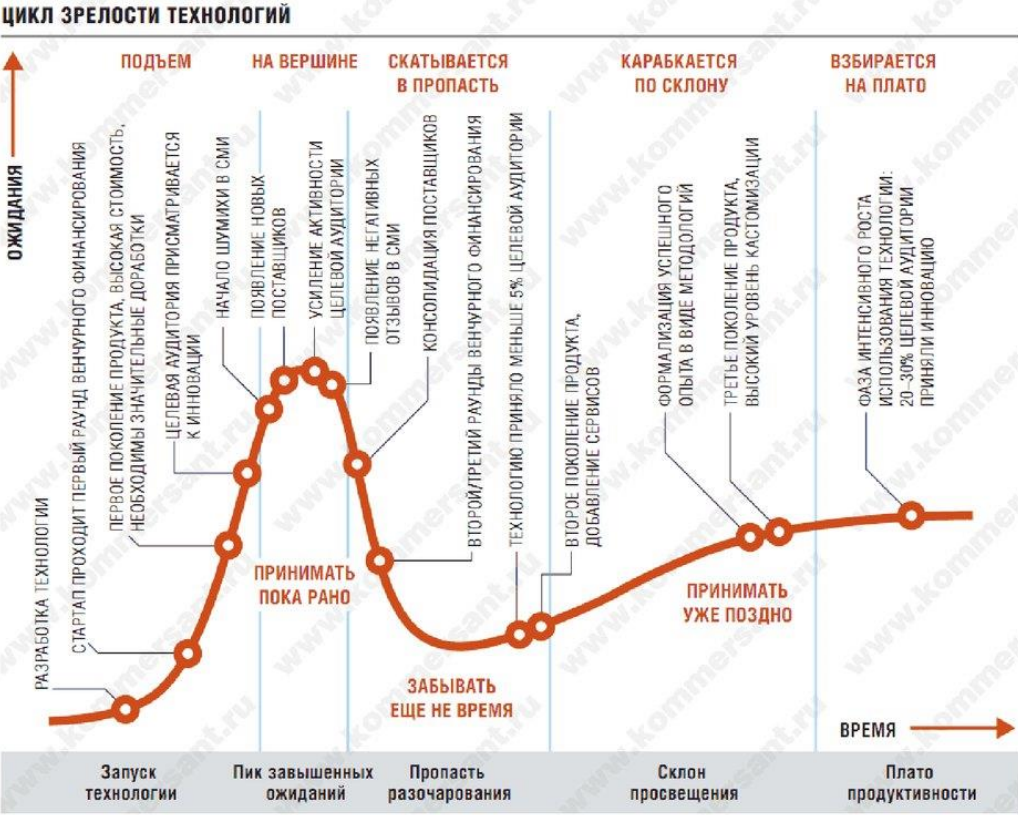
- ▶ CS231n - DL for CV (Stanford U.)
- ▶ CS224n - NLP with DL (Stanford U.)
- ▶ http://github.com/yandexdataschool/Practical_DL
- ▶ http://github.com/yandexdataschool/nlp_course
- ▶ <http://dlcourse.ai/>
- ▶ Pytorch Tutorials
- ▶ <http://kaggle.com>

Глубокое* обучение**

*глубинное
** deep learning, DL



Хайп DL



Hype Cycle for Artificial Intelligence, 2024



Deep Learning



Преимущество

- ▶ Гибкость моделирования.
- ▶ Следование данным.
- ▶ Архитектурные особенности для типа данных (а не для задачи).

Deep Learning



Преимущество

- ▶ Гибкость моделирования.
- ▶ Следование данным.
- ▶ Архитектурные особенности для типа данных (а не для задачи).

Табличные данные - (пока что?) царство градиентного бустинга

Tabular Data

columns = attributes for those observations

Player	Minutes	Points	Rebounds	Assists
A	41	20	6	5
B	30	29	7	6
C	22	7	7	2
D	26	3	3	9
E	20	19	8	0
F	9	6	14	14
G	14	22	8	3
I	22	36	0	9
J	34	8	1	3

Rows = observations

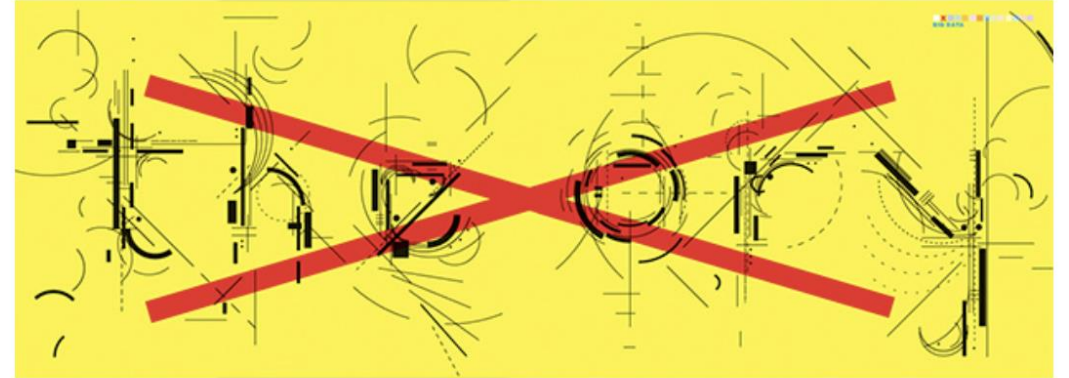
Нейросети хороши в автоматическом выделении признаков, а в табличных данных они уже выделены

Deep Learning



CHRIS ANDERSON SCIENCE 06.23.08 12:00 PM

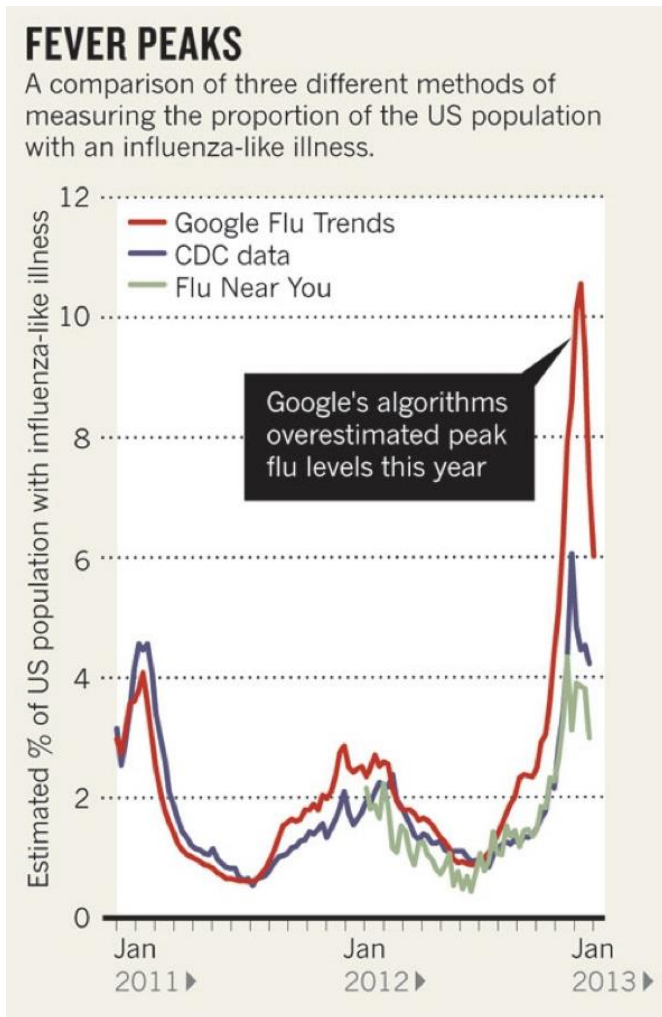
THE END OF THEORY: THE DATA DELUGE MAKES THE SCIENTIFIC METHOD OBSOLETE



* Illustration: Marian Bantjes * **"All models are wrong, but some are useful."**

<https://www.wired.com/2008/06/pb-theory/>

Deep Learning



<https://www.nature.com/news/when-google-got-flu-wrong-1.12413>

CHRIS ANDERSON SCIENCE 06.23.08 12:00 PM

THE END OF THEORY: THE DATA DELUGE MAKES THE SCIENTIFIC METHOD OBSOLETE



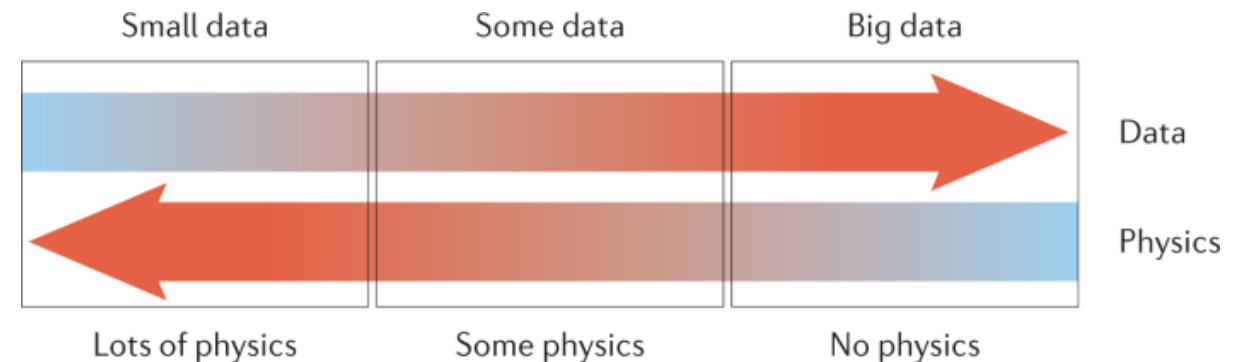
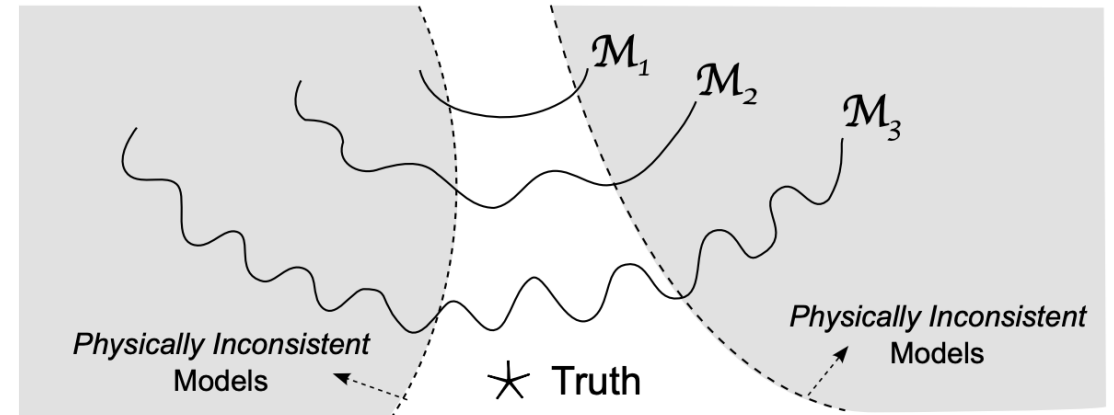
* Illustration: Marian Bantjes * **"All models are wrong, but some are useful."**

<https://www.wired.com/2008/06/pb-theory/>

Теория и глубокое обучение

- Робастность
- Экстраполяция.
- Объяснимость.
- Ясность.

Успешное применение теории в местах с маленькими выборками данных (не от хорошей жизни).

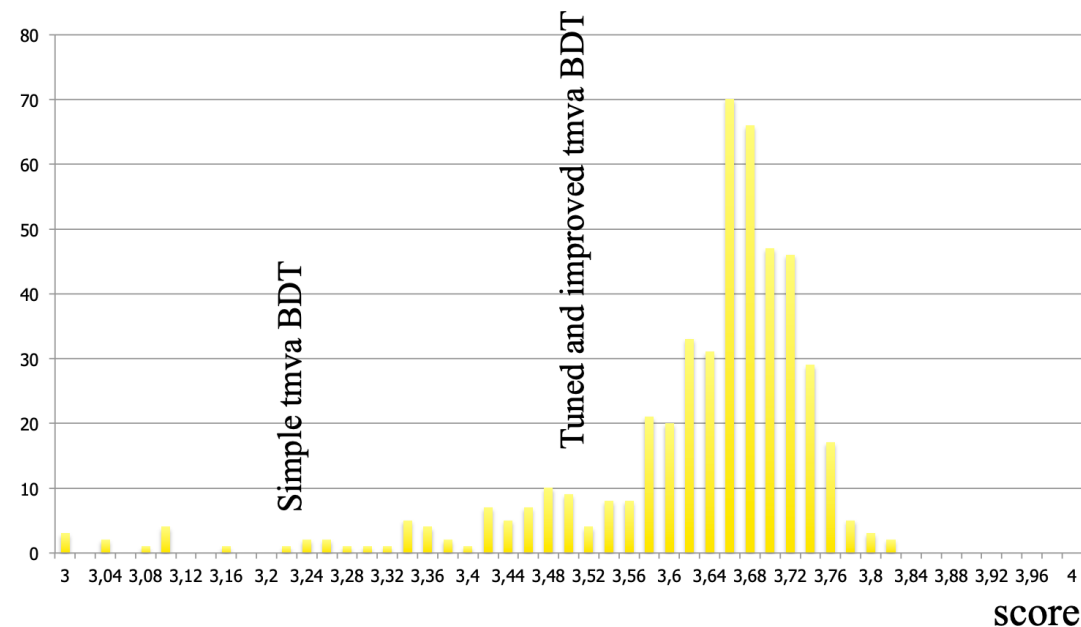


Anuj Karpatne, et al., IEEE Transactions on knowledge and data engineering 29, 10 (2017), 2318–2331

Глубокое обучение и теория

- Использование подхода «черного ящика» в приложениях.
- Быстрая адаптация к новым случаям.
- Анализ многомерных данных.

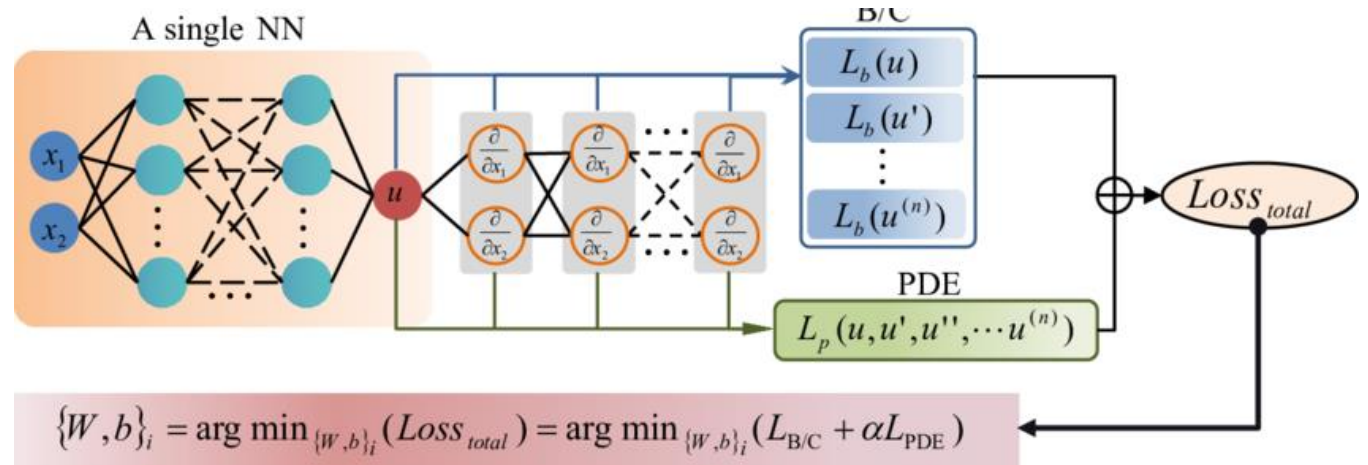
Внедрение методов глубокого обучения позволяет более детально изучить классические методы.



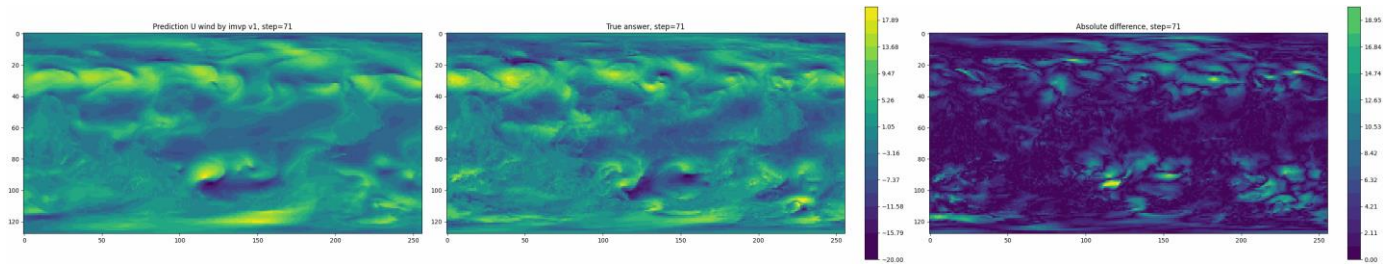
D. Roussau, Connecting the dots 2025
https://indico.physics.lbl.gov/event/149/contributions/219/attachments/211/225/tr150210_davidRoussau_CTD2015_HiggsML.pptx.pdf

Гибридный ИИ

- Выбор смеси теоретической модели и машинного обучения
- Может быть как для скрытых параметров модели, так и для наблюдаемых..



Требует много времени разработчиков. В долгосрочной перспективе - дает хорошие практические решения.



Take-home message

Глубокое обучение проникает во многие сферы человеческой деятельности.

В первую очередь, выигрыш достигается в областях с большим количеством доступных данных и слабыми моделями.

Перспектива – гибридизация глубокого обучения и классических моделей.

История глубокого обучения



Обратное распространение ошибки

1986: Rumelhart, Hinton, Williams - Backpropagation, алгоритм обратного распространения ошибки для обучения нейросетей

Вообще говоря, открывался и переоткрывался несколько раз до этого, в частности Seppo Linnainmaa (1970), Paul Werbos (1982)

- 1980: Fukushima - Неокогнитрон, первая сверточная нейросеть
- 1989: LeCun - CNN для распознавания рукописных цифр (почтовые индексы)
- 1997: Hochreiter, Schmidhuber - LSTM

Глубокое обучение

- 2006: Термин Deep Learning
- 2006: Hinton, Deep Belief Networks - успешное обучение глубоких моделей за счет послойного предобучения без учителя
- 2009: Успех в распознавании речи

Take-home message

Развитие нейронных сетей переживало взлёты и замедления.

Основой развития глубокого обучения стало накопление опыта применения классического машинного обучения и развитие математического аппарата.

К 2010-м годам было накоплено достаточно данных и компьютерных мощностей для получения революционных условий.



Машинное обучение и нейронные сети



Машинное обучение

По сути все, что мы делаем в машинном обучении - это аппроксимация и оптимизация функций:

- есть какая зависимость в реальных данных, очень сложная функция
- мы пытаемся подобрать модель (функцию), которая хорошо аппроксимирует реальную зависимость
- для этого выбираем какой-то класс моделей, обычно параметрический
- подгоняем параметры так, чтобы модель хорошо описывала данные
- то есть оптимизируем какую-то функцию ошибки или функцию качества

Нейронные сети - это очень мощный и гибкий класс моделей.

Обучение с учителем

Обучающий датасет

$$D = \{x_i, y_i\}_{i=1}^N$$

Модель - семейство функций

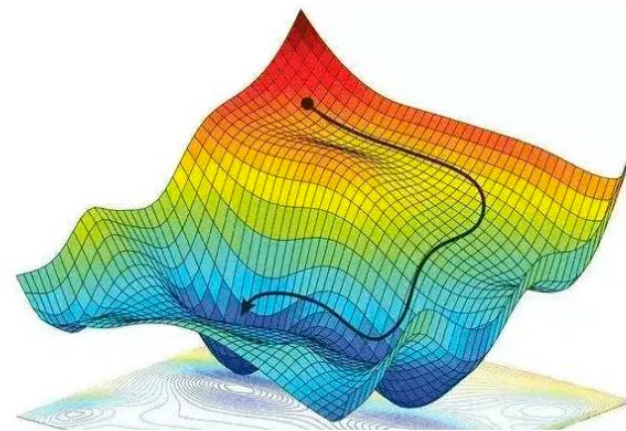
$$\hat{y} = f(x, \theta)$$

Функция потерь (Loss function)

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

Нужно настроить параметры

$$\theta^* = \underset{\theta}{\operatorname{argmin}} L(\theta)$$



$$\theta' = \theta - \alpha \frac{\partial L}{\partial \theta}$$

Обучение с учителем

Обучающий датасет

$$D = \{x_i, y_i\}_{i=1}^N$$

Модель - семейство функций

$$\hat{y} = f(x, \theta)$$

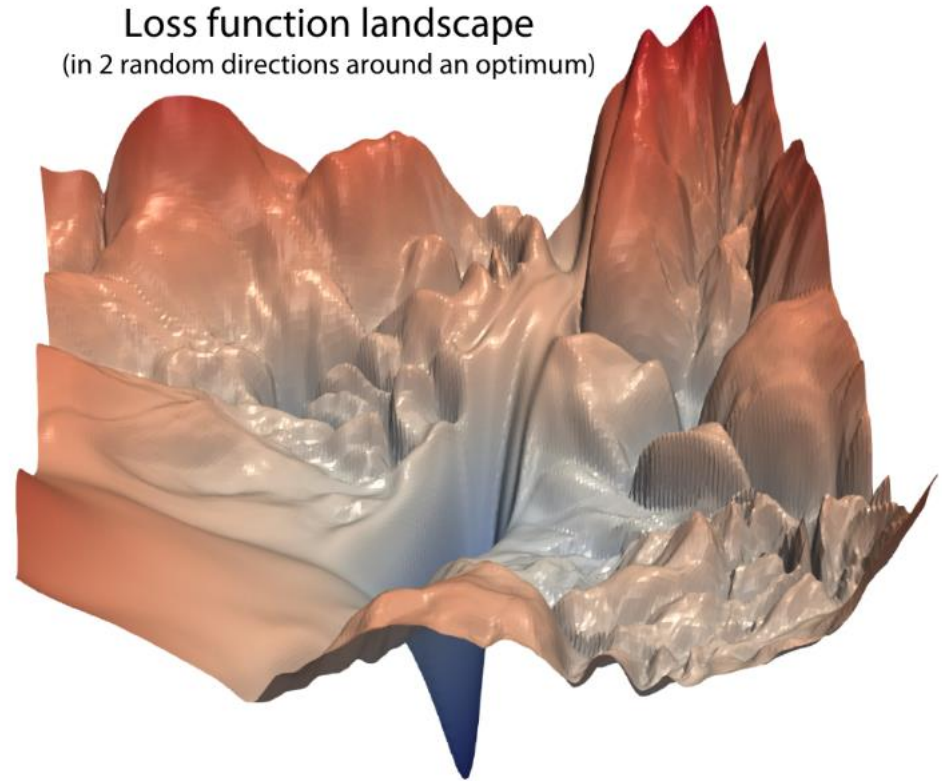
Функция потерь (Loss function)

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

Нужно настроить параметры

$$\theta^* = \operatorname{argmin}_{\theta} L(\theta)$$

Loss function landscape
(in 2 random directions around an optimum)



<https://papers.nips.cc/paper/7875-visualizing-the-loss-landscape-of-neural-nets>

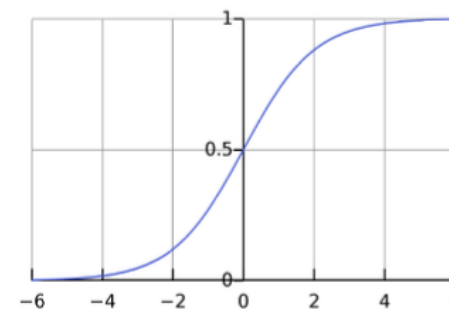
$$\theta' = \theta - \alpha \frac{\partial L}{\partial \theta}$$

Логистическая регрессия

Бинарная классификация: метки $y_i \in \{0, 1\}$

Как сделать так, чтобы
предсказания были от 0 до 1
(вероятности)?

$$P(y = 1|x, \theta) = \sigma(\theta^T x) = \frac{1}{1 + \exp(-\theta^T x)}$$



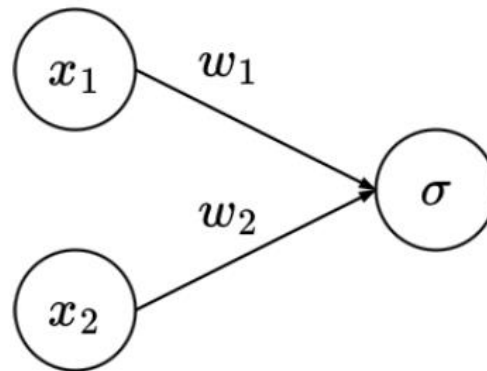
Функция потерь:

$$L(\theta) = - \sum_i (y_i \log P(y = 1|x, \theta) + (1 - y_i) \log(1 - P(y = 1|x, \theta)))$$

Получаем линейную разделяющую границу

Логистическая регрессия в графическом представлении

$$y = \sigma(w_0 + w_1x_1 + w_2x_2)$$

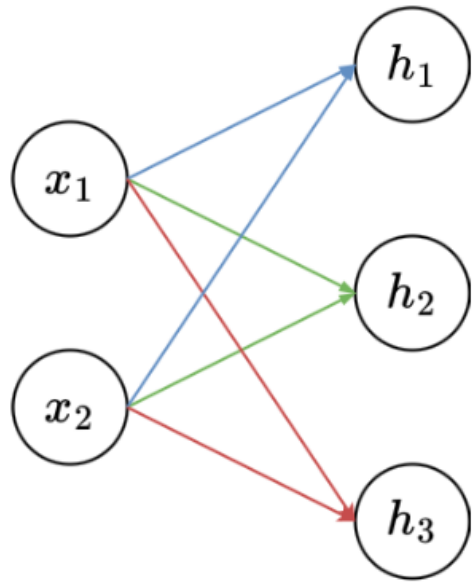


Логистическая регрессия = 1 нейрон

Вход -> линейная операция -> нелинейность -> выход

Несколько регрессий в одной форме

Составим три логистических регрессии с одними входами



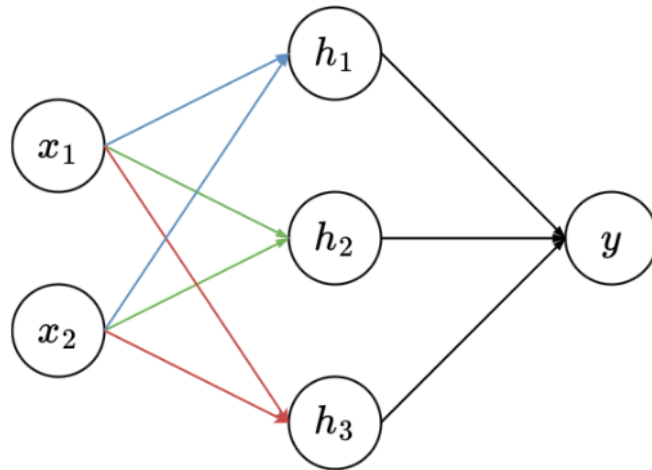
$$h_1 = \sigma(b_1 + w_{11}x_1 + w_{12}x_2)$$

$$h_2 = \sigma(b_2 + w_{21}x_1 + w_{22}x_2)$$

$$h_3 = \sigma(b_3 + w_{31}x_1 + w_{32}x_2)$$

Логистическая регрессия в графическом представлении

Объединим выходы трёх логистических регрессий в единый выход



$$h_1 = \sigma(b_1 + w_{11}x_1 + w_{12}x_2)$$

$$h_2 = \sigma(b_2 + w_{21}x_1 + w_{22}x_2)$$

$$h_3 = \sigma(b_3 + w_{31}x_1 + w_{32}x_2)$$

$$y = \sigma(b' + w'_{11}h_1 + w'_{12}h_2 + w'_{13}h_3)$$

Матричное представление

$$h_1 = \sigma(b_1 + w_{11}x_1 + w_{12}x_2)$$

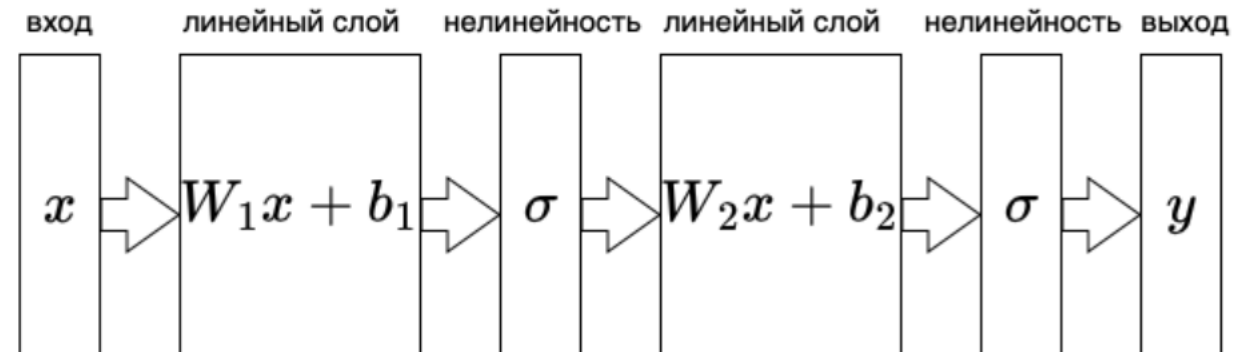
$$h_2 = \sigma(b_2 + w_{21}x_1 + w_{22}x_2)$$

$$h_3 = \sigma(b_3 + w_{31}x_1 + w_{32}x_2)$$

$$h = \sigma(Wx + b)$$

$$y = \sigma(W'h + b')$$

$$y = \sigma(b' + w'_{11}h_1 + w'_{12}h_2 + w'_{13}h_3)$$

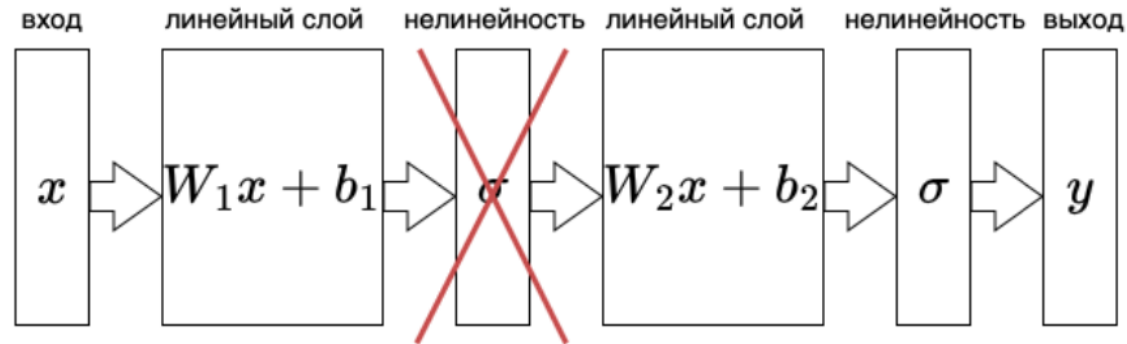


$$y = \sigma(W_2\sigma(W_1x + b_1) + b_2)$$

Логистическая регрессия в графическом представлении

Что если уберем нелинейность?

Суперпозиция линейных функций - снова линейная функция.



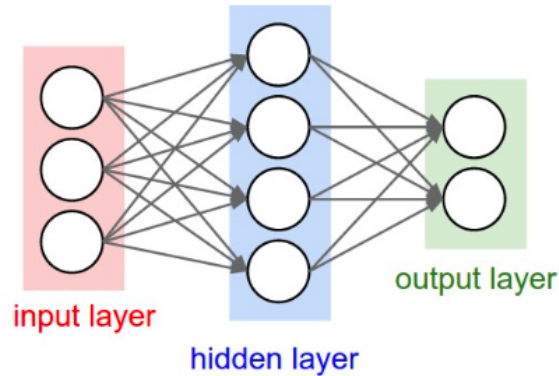
$$y = \sigma(W_2(W_1x + b_1) + b_2) = \sigma(W_2W_1x + W_2b_1 + b_2) = \sigma(W'x + b')$$

Но! Не все нелинейности одинаково полезны!

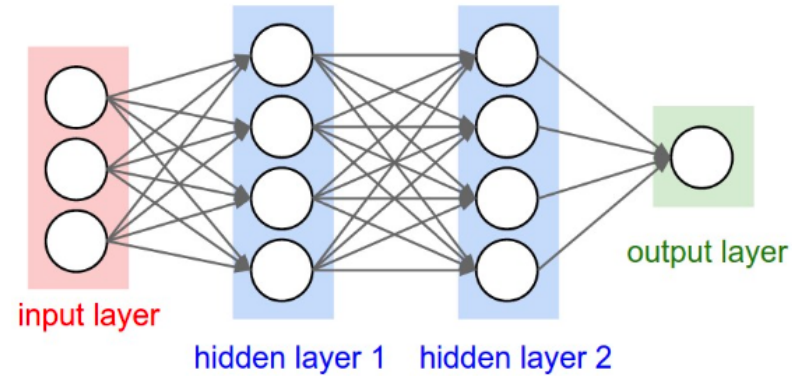
Multilayer Perceptron (MLP)

Многослойный перцептрон

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



$$y = \sigma(W_3\sigma(W_2\sigma(W_1x + b_1) + b_2) + b_3)$$



W - веса (weights), b - смещение (bias), σ - функция активации

Организация в слои - можно представить в виде матричных операций и легко параллелизовать

Развитие глубокого обучения

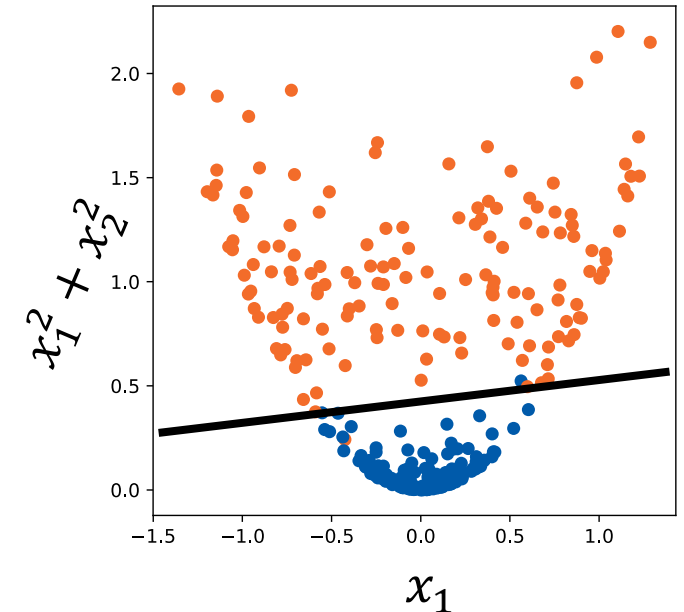
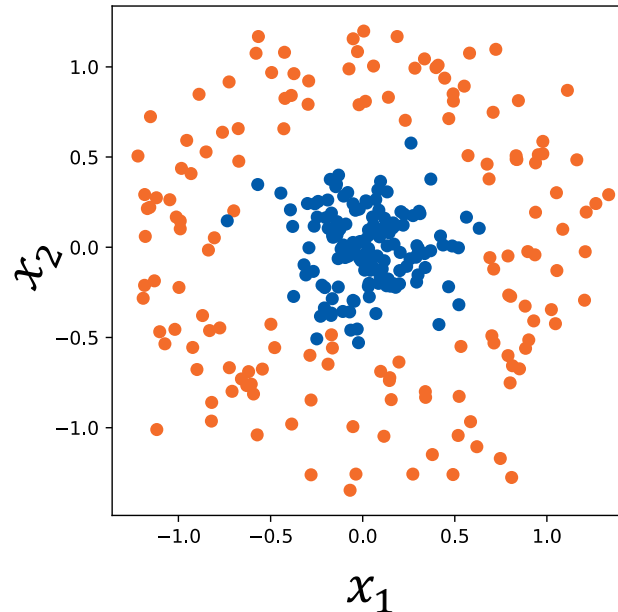
- Вычислительные мощности (GPUs, TPUs)
- Огромное количество данных
- Улучшение алгоритмов и понимания того, как обучать глубокие сети
- Open source инструменты и модели
- Интерес бизнеса, больших корпораций

От линейной модели к нейронной сети



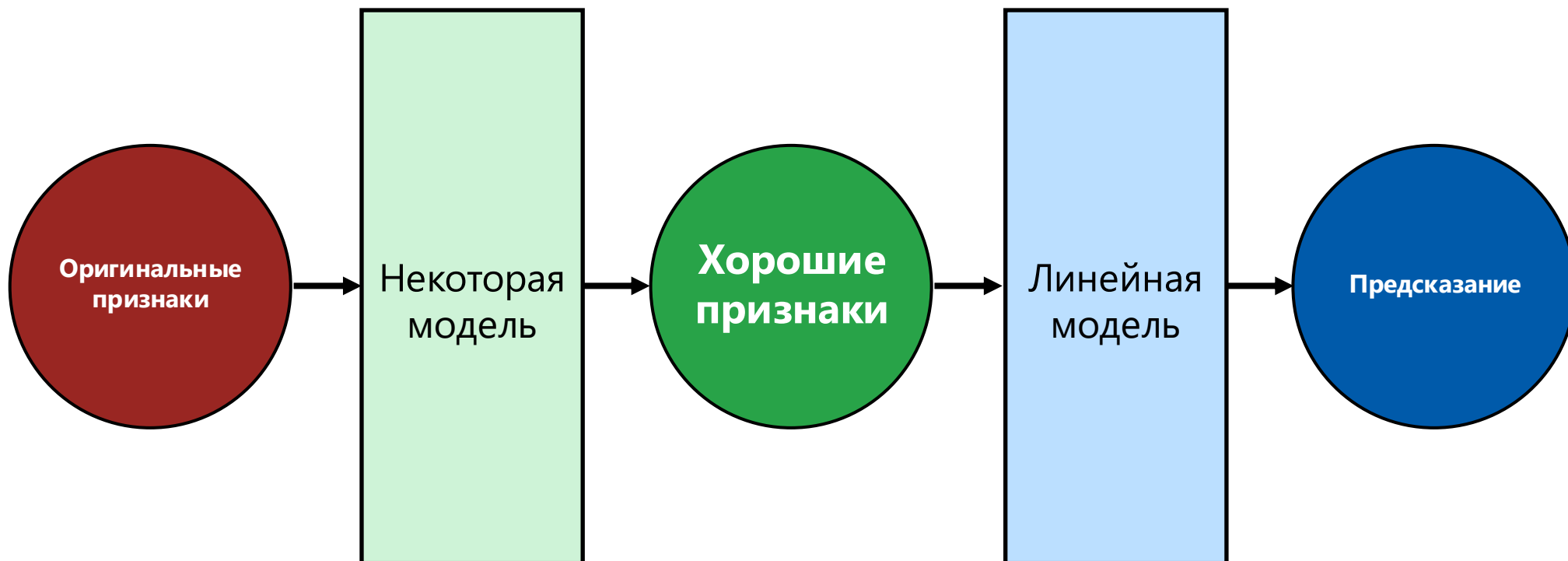
Линейные модели в извлечение признаков

- ▶ Для линейных моделей мы **придумываем** новые признаки, чтобы сделать модель более мощной
- ▶ Поиск хороших функций (так называемая **feature engineering**) - задача крайне нетривиальная.



**Автоматизация конструирования признаков –
основная сила DL.**

Идея: стэкинг моделей



- ▶ Добавить ещё одну модель
- ▶ Всё обучим одновременно
 - Если модели **дифференцируемые**, можем использовать градиентный спуск

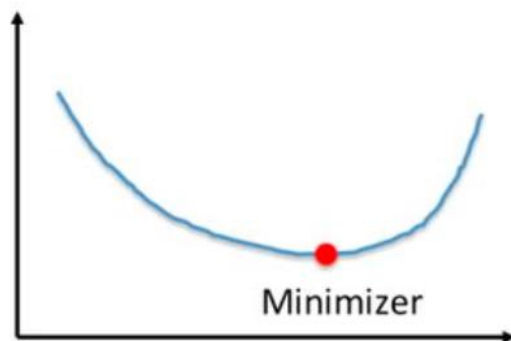
Замечание: такое суммирование моделей, вероятно, делает задачу невыпуклой
⇒ **нет гарантий сходимости**

Convex vs non-convex optimization

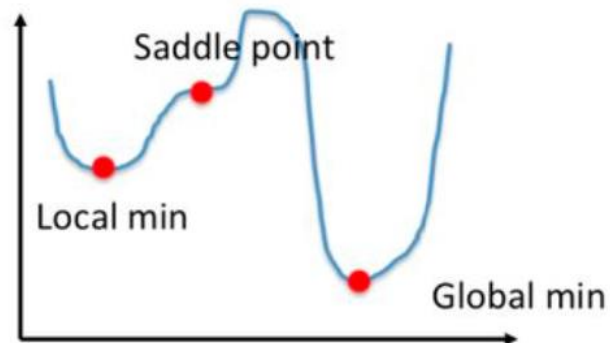
Выпуклая оптимизация

Невыпуклая оптимизация

Convex



Non-Convex

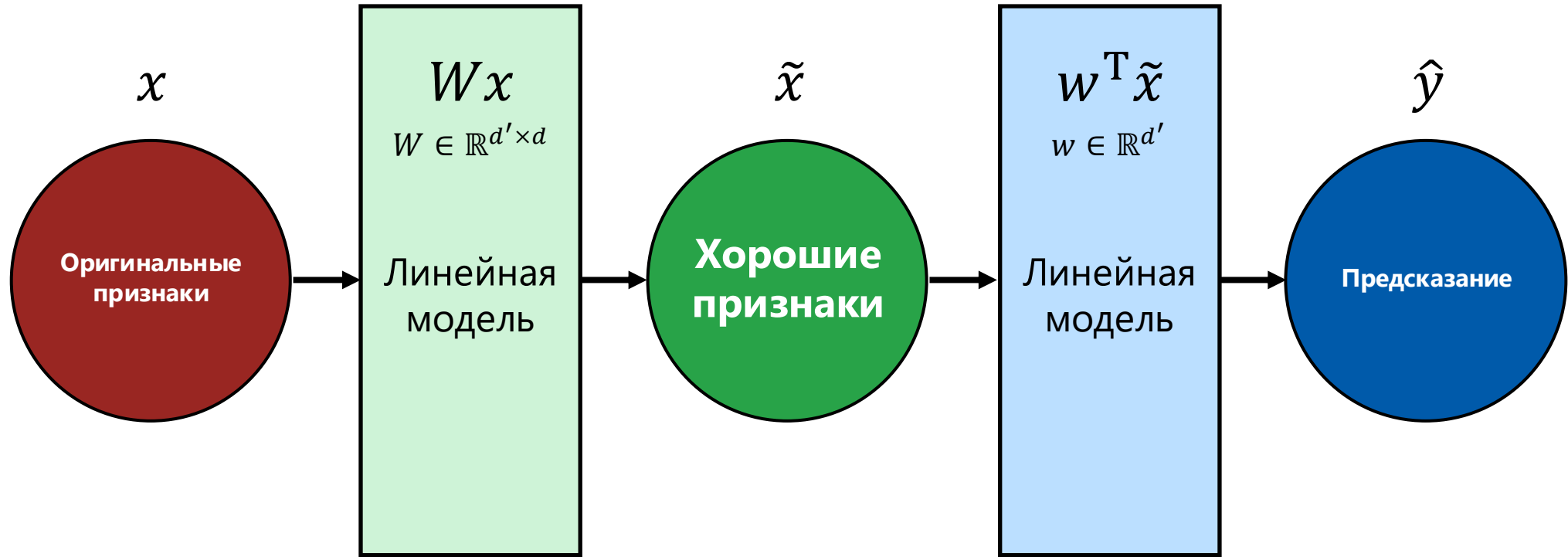


<https://ru.pinterest.com/pin/convex-nonconvex-functions--672232681861372201/>

Множество минимумов, неявные точки перегиба

Замечание: такое суммирование моделей, вероятно, делает задачу невыпуклой
⇒ **нет гарантий сходимости**

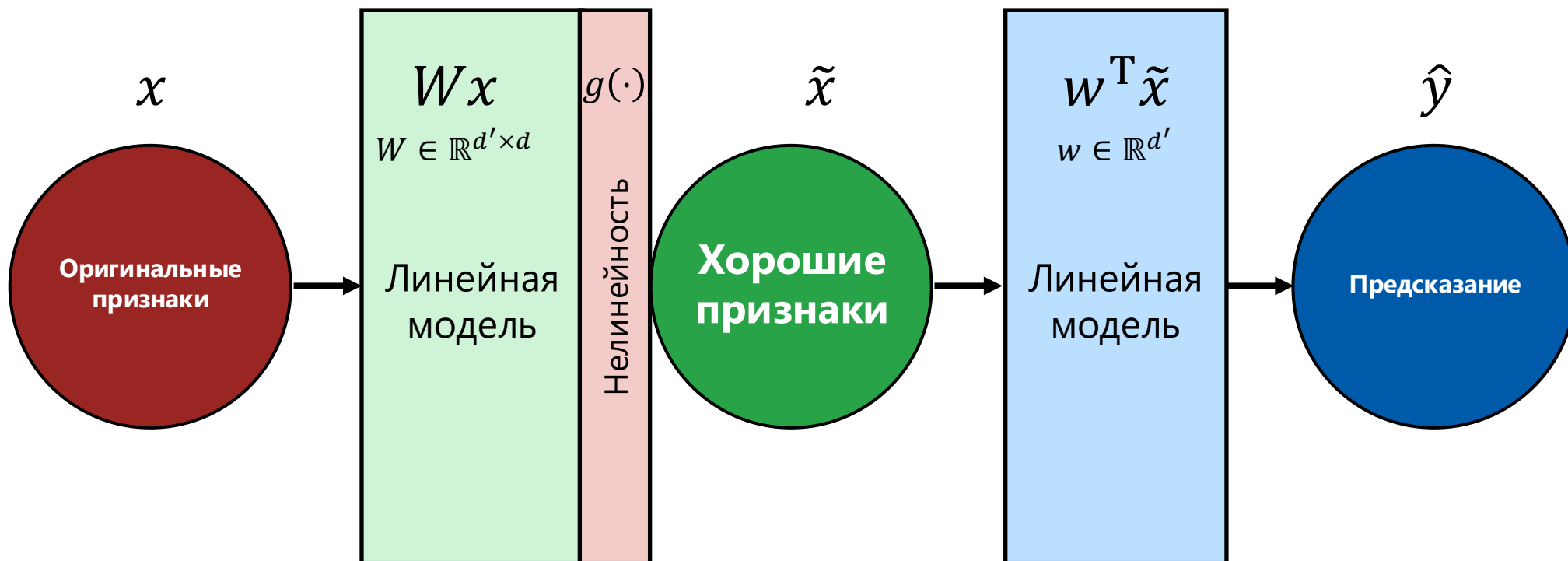
Какие модели идут в стэк?



$$\hat{y} = w^T \tilde{x} = w^T (Wx) = (w^T W)x = w'^T x$$

– Всё становится линейной моделью

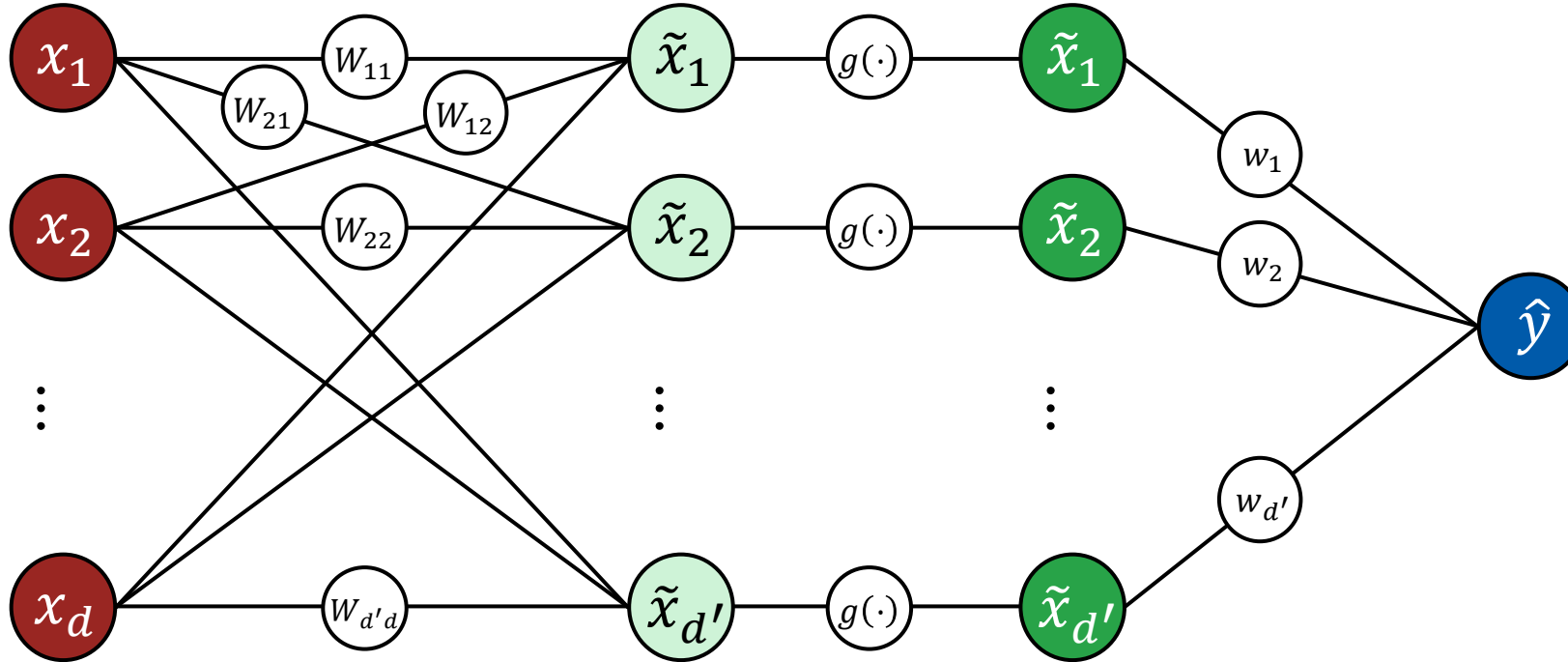
Введём нелинейность



$$\hat{y} = w^T \tilde{x} = w^T g(Wx)$$

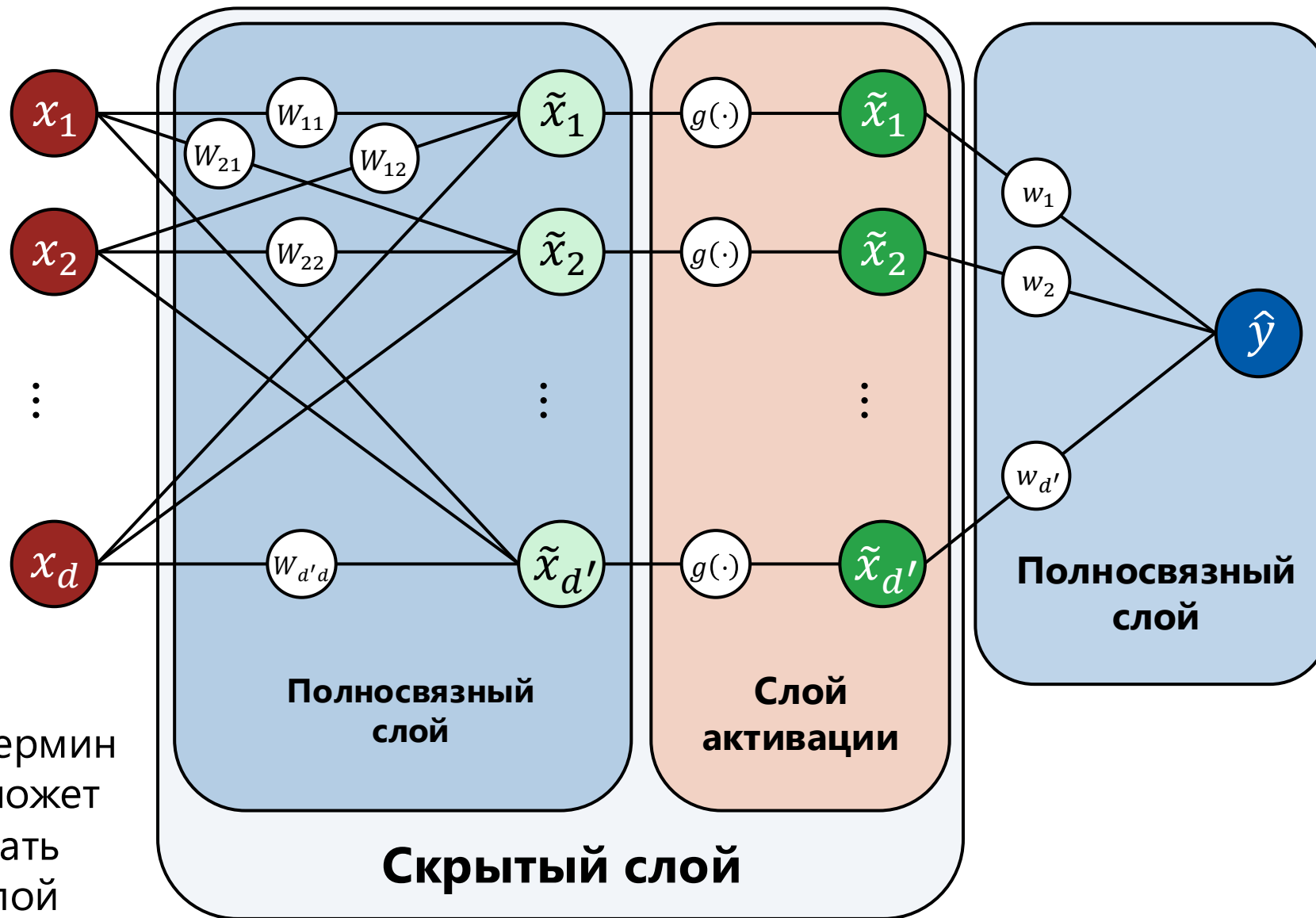
$g(\cdot)$ – некоторая **нелинейная** скалярная функция (поэлементная)

Детализированный вид



$$\hat{y} = w^T \tilde{x} = w^T g(Wx) = \sum_j \left[w_j g \left(\sum_i W_{ji} x_i \right) \right]$$

Терминология

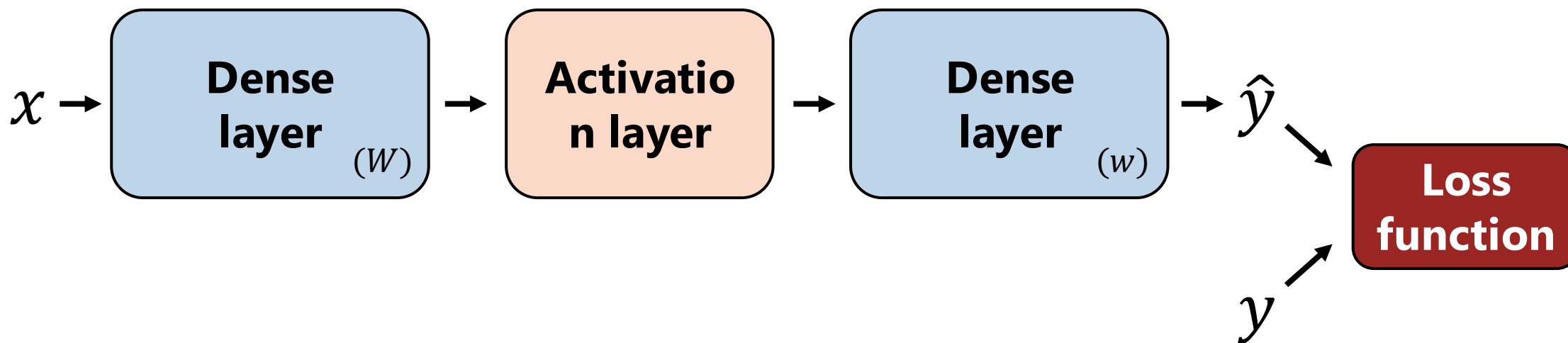


Примечание: термин
«активация» может
также означать
выходной слой

Обратное распространение ошибки



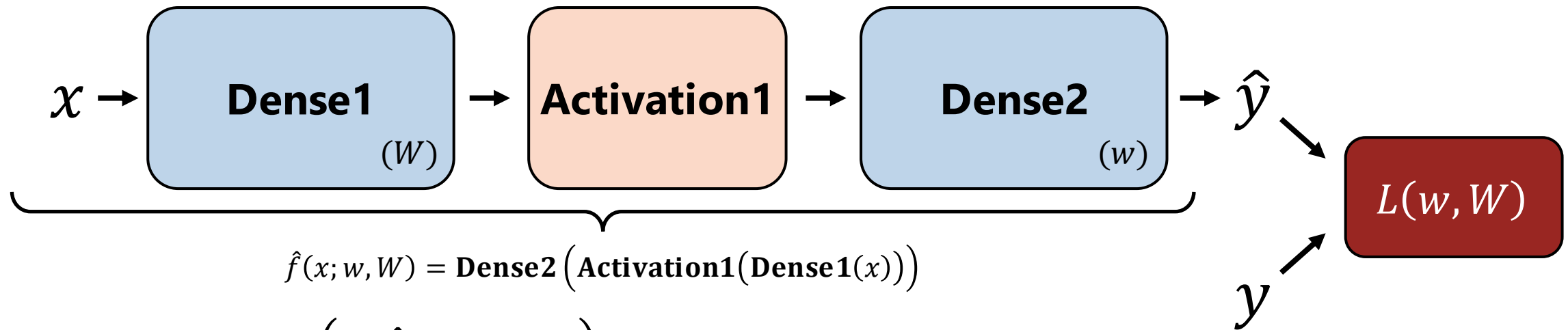
Функция потерь



- ▶ Например: квадрат ошибки

$$L = \frac{1}{N} \sum_{i=1 \dots N} \left(y_i - w^T g(W x_i) \right)^2$$

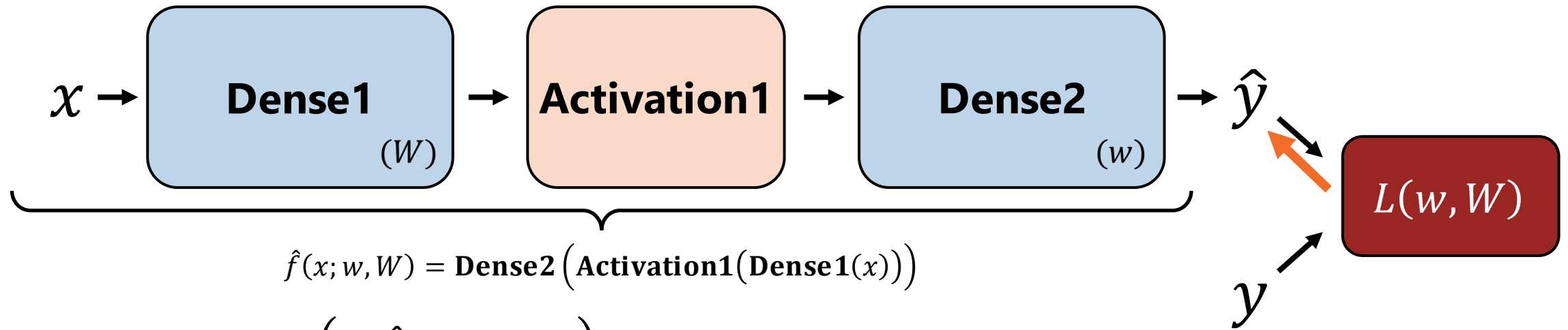
Производные



$$L(w, W) \equiv L(y, \hat{f}(x; w, W))$$

$$\frac{\partial L}{\partial W} = \frac{\partial L(y, \hat{f})}{\partial \hat{f}} \cdot \frac{\partial \hat{f}}{\partial W} =$$

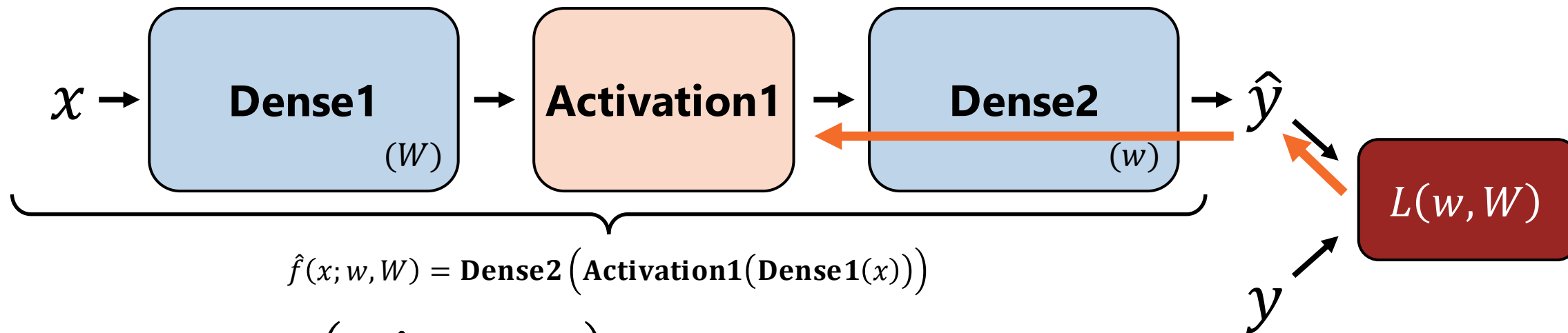
Производные



$$L(w, W) \equiv L(y, \hat{f}(x; w, W))$$

$$\frac{\partial L}{\partial W} = \frac{\partial L(y, \hat{f})}{\partial \hat{f}} \cdot \frac{\partial \hat{f}}{\partial W} = \frac{\partial L(y, \hat{f})}{\partial \hat{f}} \cdot$$

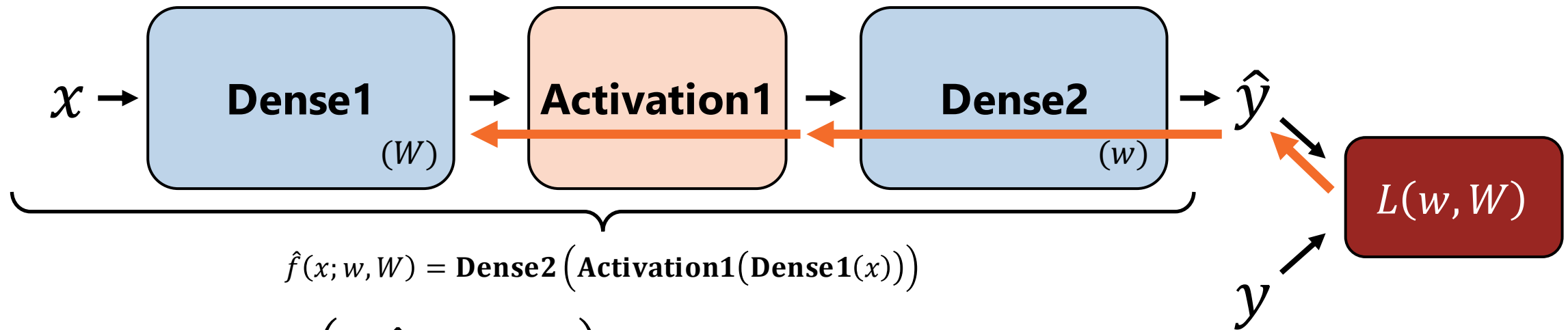
Снова производные



$$L(w, W) \equiv L(y, \hat{f}(x; w, W))$$

$$\frac{\partial L}{\partial W} = \frac{\partial L(y, \hat{f})}{\partial \hat{f}} \cdot \frac{\partial \hat{f}}{\partial W} = \frac{\partial L(y, \hat{f})}{\partial \hat{f}} \cdot \frac{\partial \text{Dense2}}{\partial \text{Activation1}}.$$

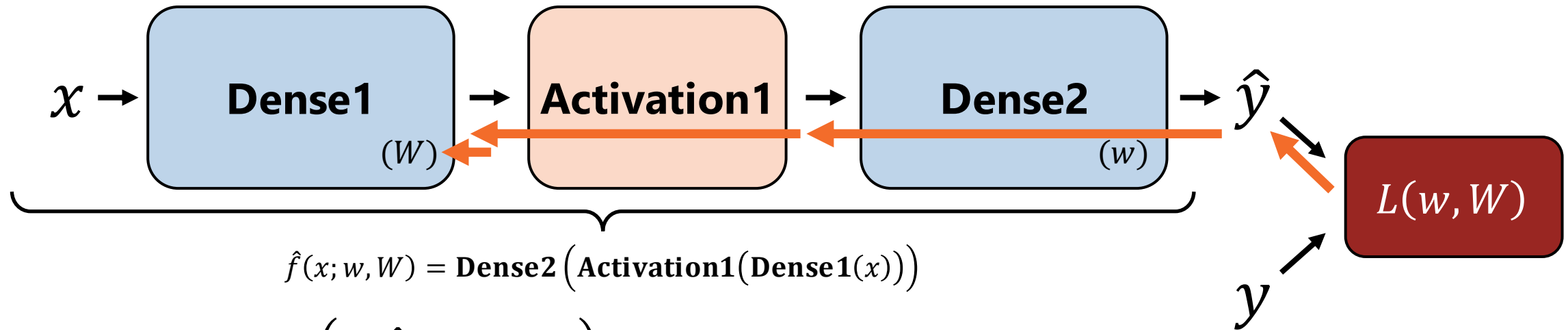
Calculating derivatives



$$L(w, W) \equiv L(y, \hat{f}(x; w, W))$$

$$\frac{\partial L}{\partial W} = \frac{\partial L(y, \hat{f})}{\partial \hat{f}} \cdot \frac{\partial \hat{f}}{\partial W} = \frac{\partial L(y, \hat{f})}{\partial \hat{f}} \cdot \frac{\partial \text{Dense2}}{\partial \text{Activation1}} \cdot \frac{\partial \text{Activation1}}{\partial \text{Dense1}} \cdot$$

Calculating derivatives

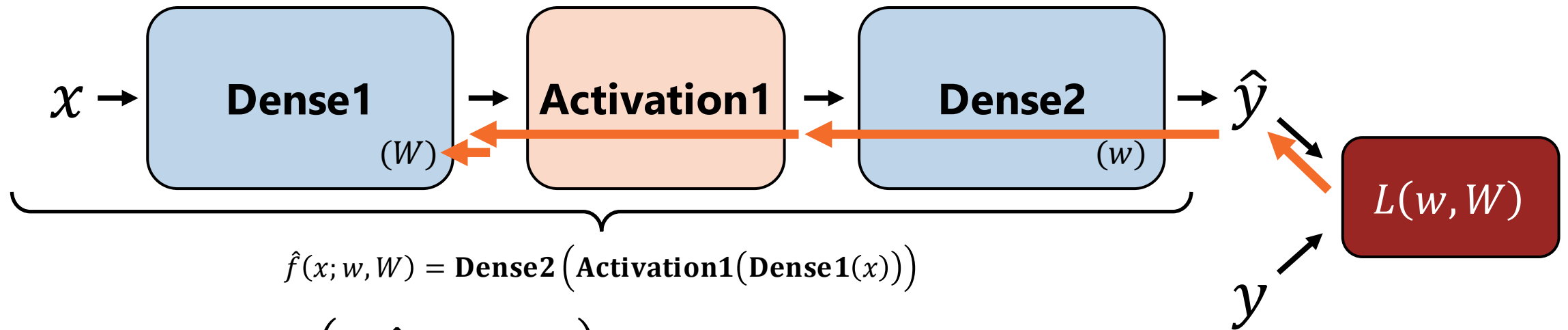


$$L(w, W) \equiv L(y, \hat{f}(x; w, W))$$

$$\frac{\partial L}{\partial W} = \frac{\partial L(y, \hat{f})}{\partial \hat{f}} \cdot \frac{\partial \hat{f}}{\partial W} = \frac{\partial L(y, \hat{f})}{\partial \hat{f}} \cdot \frac{\partial \text{Dense2}}{\partial \text{Activation1}} \cdot \frac{\partial \text{Activation1}}{\partial \text{Dense1}} \cdot \frac{\partial \text{Dense1}}{\partial W}$$

- Backpropagation algorithm \approx applying the chain rule
 - The actual algorithm states how to do it efficiently

Calculating derivatives



$$L(w, W) \equiv L(y, \hat{f}(x; w, W))$$

$$\frac{\partial L}{\partial W} = \frac{\partial L(y, \hat{f})}{\partial \hat{f}} \cdot \frac{\partial \hat{f}}{\partial W} = \underbrace{\frac{\partial L(y, \hat{f})}{\partial \hat{f}}}_{\text{scalar}} \cdot \underbrace{\frac{\partial \text{Dense2}}{\partial \text{Activation1}}}_{d' \text{-vector}} \cdot \underbrace{\frac{\partial \text{Activation1}}{\partial \text{Dense1}}}_{d' \times d' \text{-matrix (scalar} \cdot \text{unit matrix)}} \cdot \underbrace{\frac{\partial \text{Dense1}}{\partial W}}_{d' \times d' \times d \text{-tensor} = d' \times d \text{-matrix}}$$

► Обратное распространение ошибки \approx цепное правило

– Алгоритм показывает как делать это эффективно