

Машинное обучение: нейронные сети и глубокое обучение (введение)

Эмили Драль

MADE academy, Москва 2020

Алгоритмы машинного обучения

1. Обучение с учителем: линейные модели
2. Обучение с учителем: ансамбли моделей
3. **Обучение с учителем: нейросетевые модели**
4. Обучение без учителя: обзор методов
5. Обучение с подкреплением
6. (optional) Рекомендательные системы

Нейронные сети и глубокое обучение

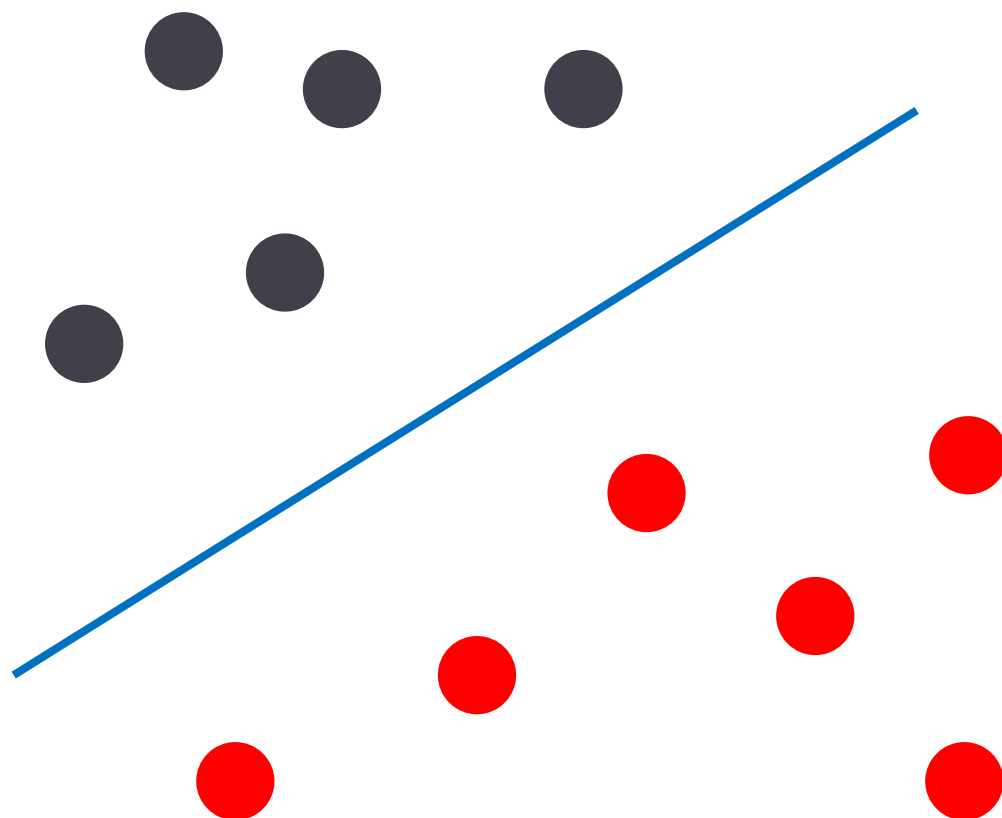
1. Связь с линейными моделями
2. Модель нейронной сети
3. Обучение модели
4. Оптимизация модели
5. Популярные приложения
6. Практические рекомендации

Связь с линейными моделями

Линейная классификация

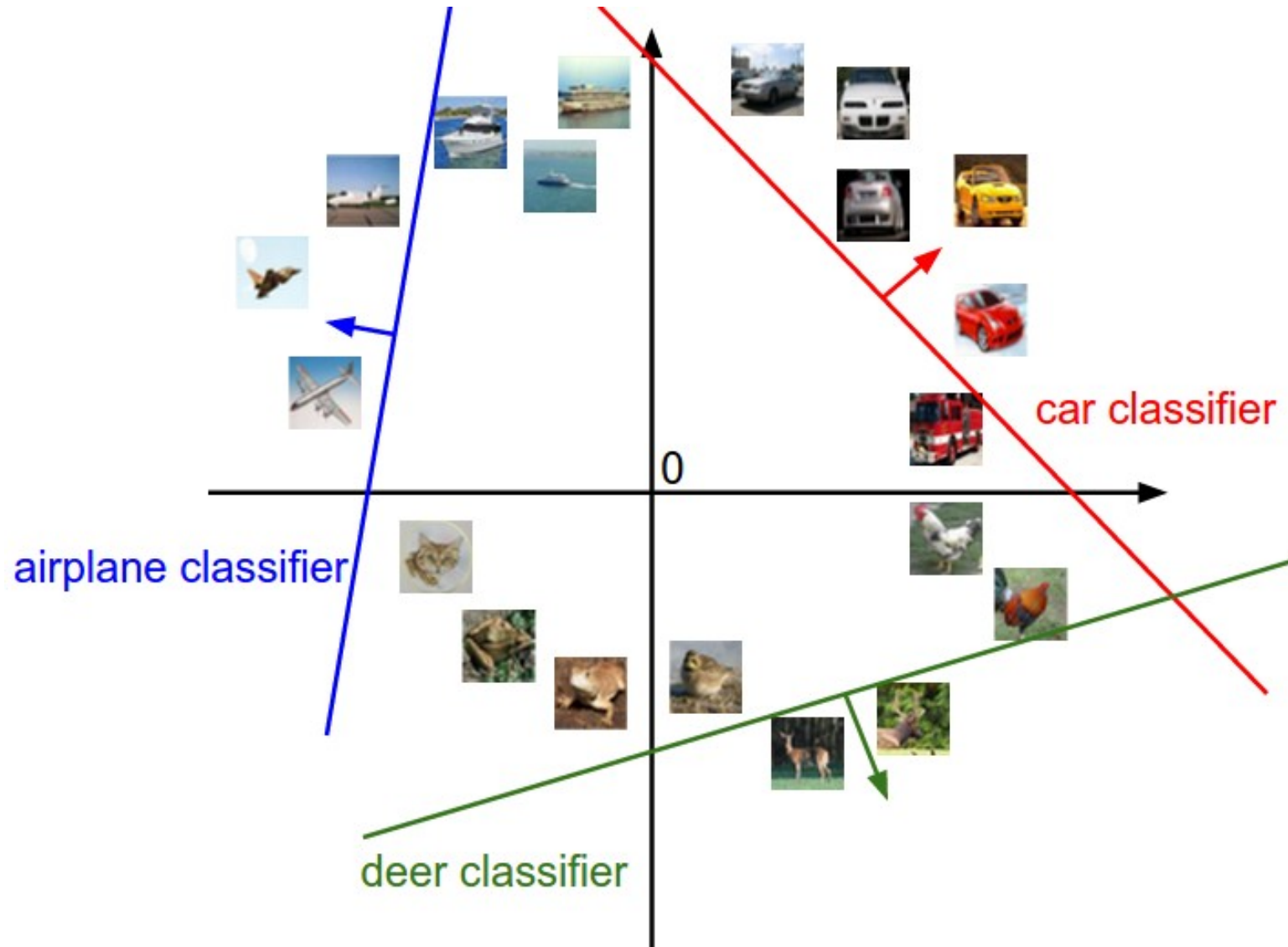
$$a(x) = \text{sign}(\langle w, x \rangle + w_0)$$

Связь с
линейными
моделями



Линейная классификация

Связь с
линейными
моделями



Как компьютер видит изображение?

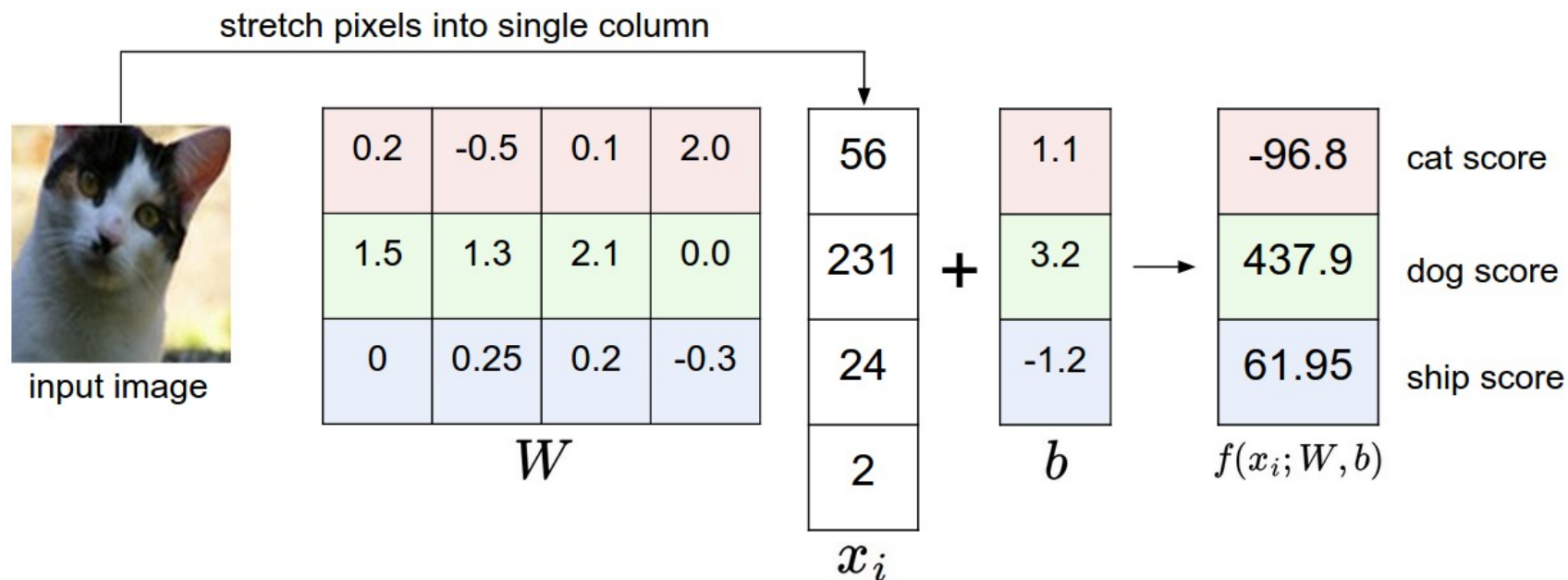
Связь с
линейными
моделями



```
08 02 22 97 38 15 00 40 00 75 04 05 07 78 52 12 50 77 91 08
49 49 99 40 17 81 18 57 60 87 17 40 98 43 69 48 04 56 62 00
81 49 31 73 55 79 14 29 93 71 40 67 53 88 30 03 49 13 36 65
52 70 95 23 04 60 11 42 69 24 68 56 01 32 56 71 37 02 36 91
22 31 16 71 51 67 63 89 41 92 36 54 22 40 40 28 66 33 13 80
24 47 32 60 99 03 45 02 44 75 33 53 78 36 84 20 35 17 12 50
32 98 81 28 64 23 67 10 26 38 40 67 59 54 70 66 18 38 64 70
67 26 20 68 02 62 12 20 95 63 94 39 63 08 40 91 66 49 94 21
24 55 58 05 66 73 99 26 97 17 78 78 96 83 14 88 34 89 63 72
21 36 23 09 75 00 76 44 20 45 35 14 00 61 33 97 34 31 33 95
78 17 53 28 22 75 31 67 15 94 03 80 04 62 16 14 09 53 56 92
16 39 05 42 96 35 31 47 55 58 88 24 00 17 54 24 36 29 85 57
86 56 00 48 35 71 89 07 05 44 44 37 44 60 21 58 51 54 17 58
19 80 81 68 05 94 47 69 28 73 92 13 86 52 17 77 04 89 55 40
04 52 08 83 97 35 99 16 07 97 57 32 16 26 26 79 33 27 98 66
88 36 68 87 57 62 20 72 03 46 33 67 46 55 12 32 63 93 53 69
04 42 16 73 38 25 39 11 24 94 72 18 08 46 29 32 40 62 76 36
20 69 36 41 72 30 23 88 34 62 99 69 82 67 59 85 74 04 36 16
20 73 35 29 78 31 90 01 74 31 49 71 48 86 81 16 23 57 05 54
01 70 54 71 83 51 54 69 16 92 33 48 61 43 52 01 89 19 67 48
```

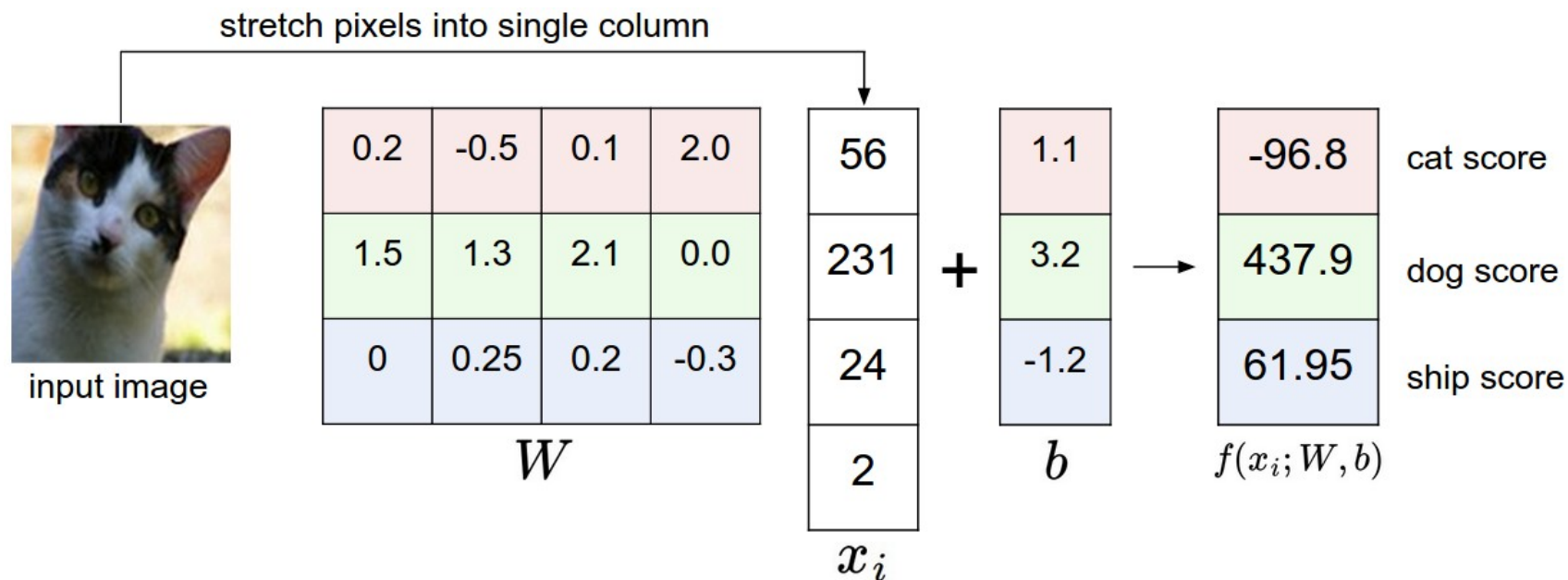
Применим линейную модель для классификации изображения

Связь с
линейными
моделями



Применим линейную модель для классификации изображения

Связь с линейными моделями



Веса линейных моделей можно интерпретировать.

Что произойдет если мы возьмем вектор весов и применим операцию `resize` для получения матрицы с размерами исходного изображения?

Шаблоны

Связь с
линейными
моделями



Связь с
линейными
моделями

В чем проблема такого
подхода?

Давайте подумаем, насколько данный подход можно
обобщить.

Потенциальные проблемы: фон

Связь с
линейными
моделями



Потенциальные проблемы: смещения, повороты

Связь с
линейными
моделями



Потенциальные проблемы: много объектов на одном изображении

Связь с
линейными
моделями



Потенциальные проблемы: объект частично скрыт

Связь с
линейными
моделями



Задачи
машинного
обучения

Идея: усложним модель?

Задачи машинного обучения

Идея: усложним модель?

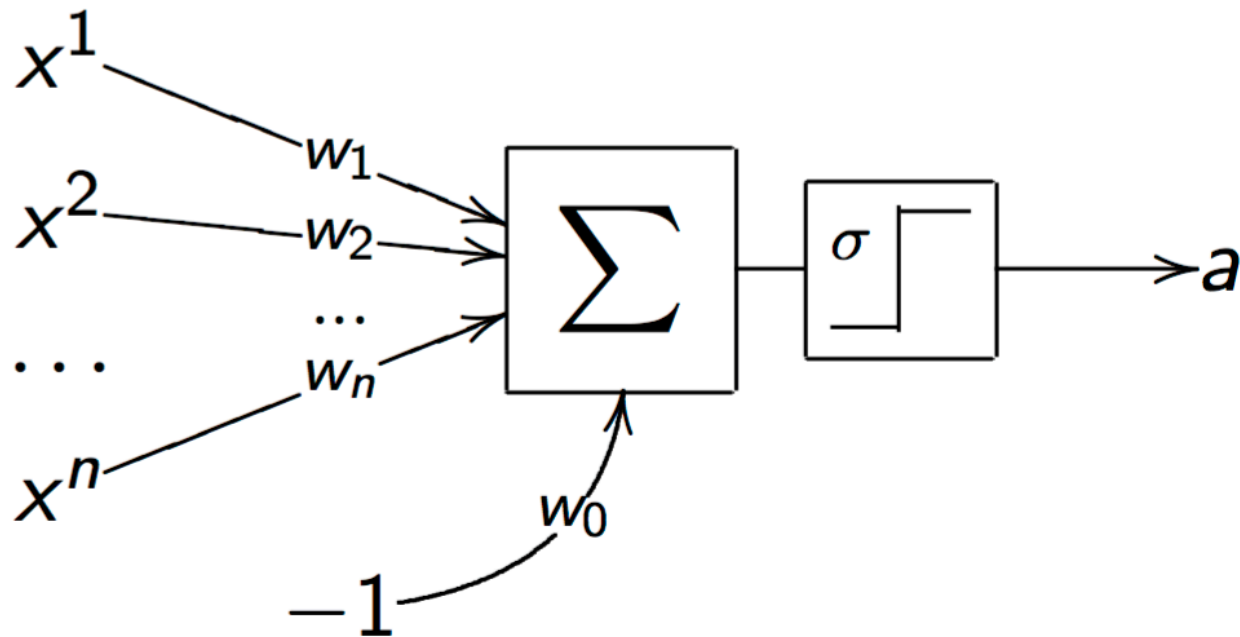
Давайте построим суперпозицию линейных моделей:

- Будем строить линейную модель над линейными моделями
- Это можно делать много раз (не только 2!)
- Надо будет добавить нелинейность, ведь суперпозиция линейных моделей оставляет нас в пространстве линейных моделей

Модель нейронной сети

Модель нейронной сети

Линейная модель нейрона



Чтобы получить линейный классификатор, достаточно
взять $\sigma = \text{sign}(z)$

Модель нейронной сети

Усложняем модель

1. Строим линейную модель над линейной моделью

$a(x) = W2 \cdot (W1 \cdot x + b1) + b2 = W \cdot x + b$ – снова
линейная модель

Модель нейронной сети

Усложняем модель

1. Строим линейную модель над линейной моделью

$a(x) = W_2 \cdot (W_1 \cdot x + b_1) + b_2 = W \cdot x + b$ — снова линейная модель

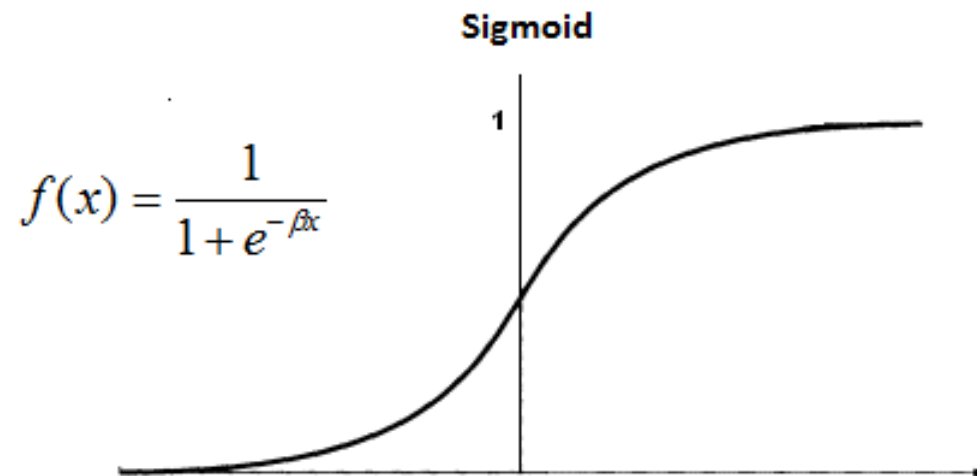
2. Добавляем нелинейность

$a(x) = W_2 \cdot \sigma(W_1 \cdot x + b_1) + b_2 \neq W \cdot x + b$ — это уже не линейная модель, значит у нас получилось усложнить

Модель нейронной сети

Усложняем модель

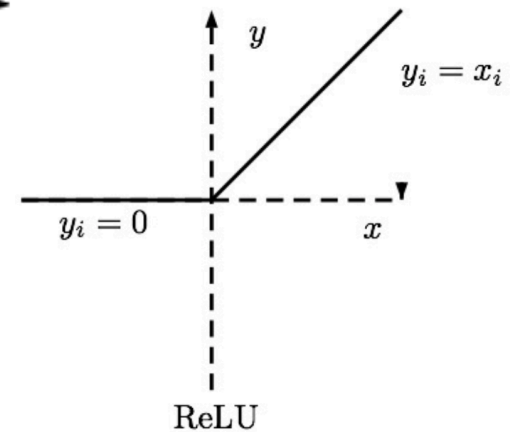
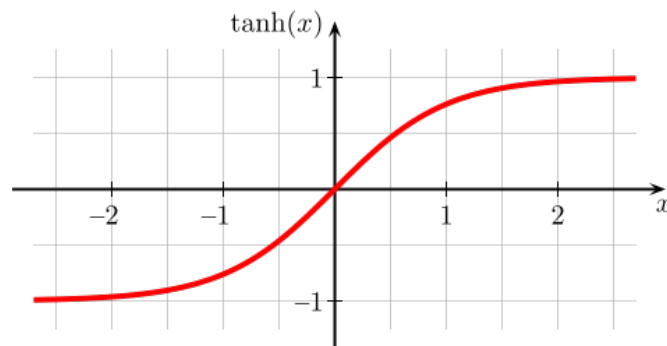
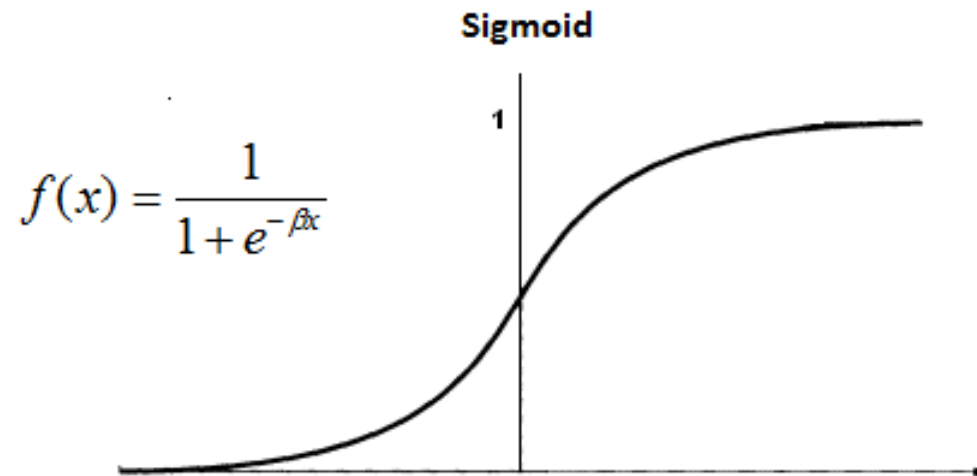
Как ввести нелинейность?



Модель нейронной сети

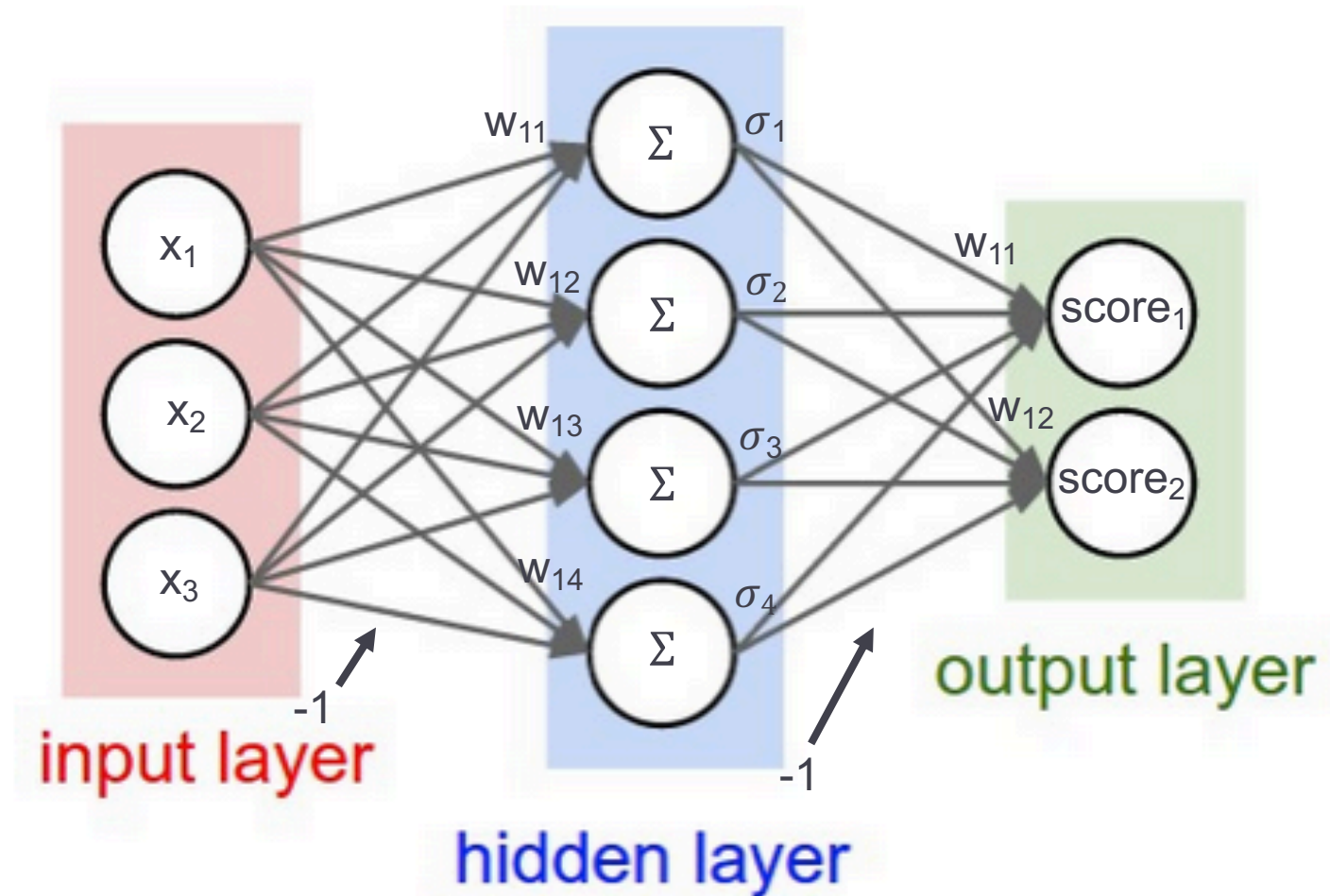
Усложняем модель

Как ввести нелинейность?



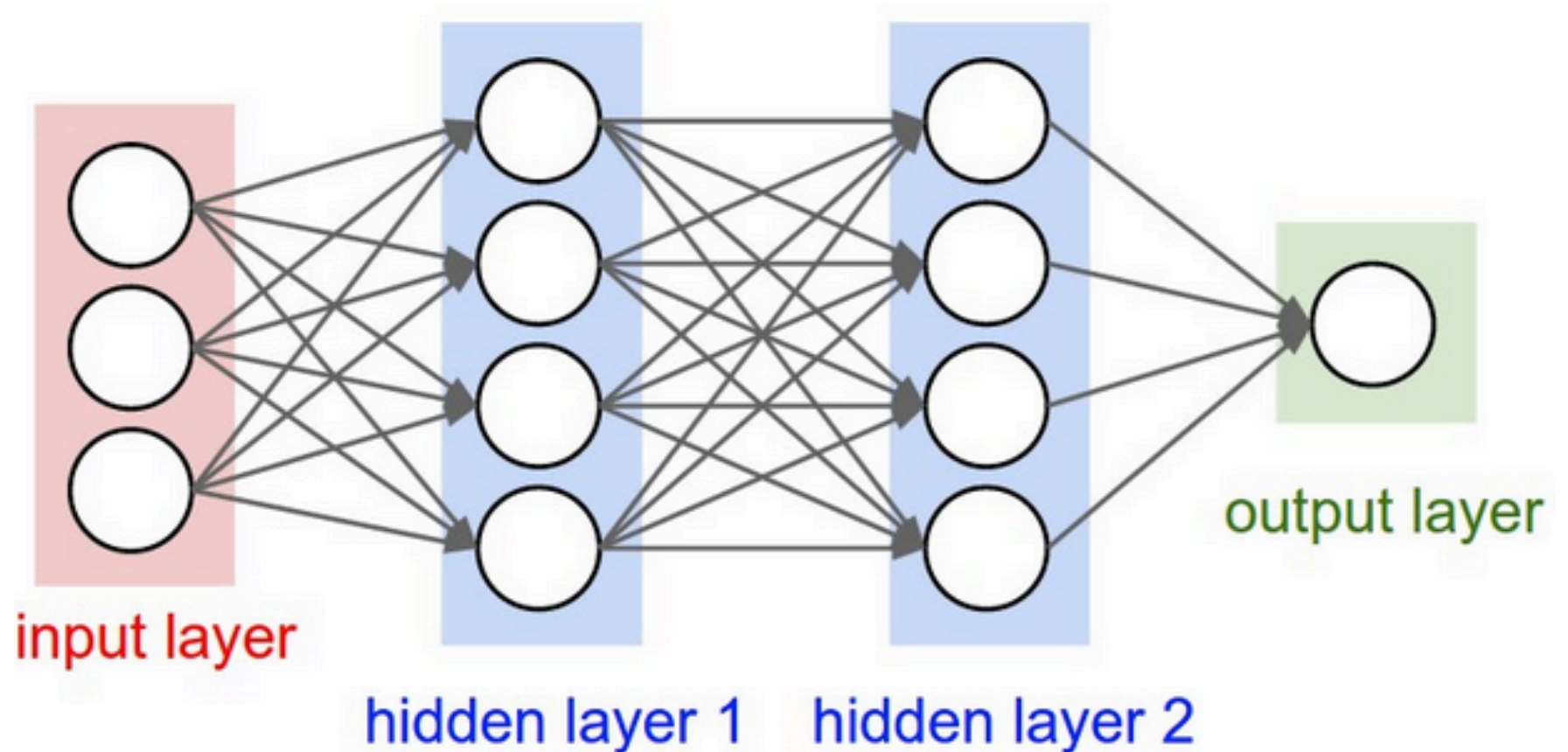
Модель нейронной сети

Собираем полносвязную нейронную сеть



Модель нейронной сети

Собираем полносвязную нейронную сеть



Обучение модели

Обучение
модели

Задача оптимизации для линейно модели

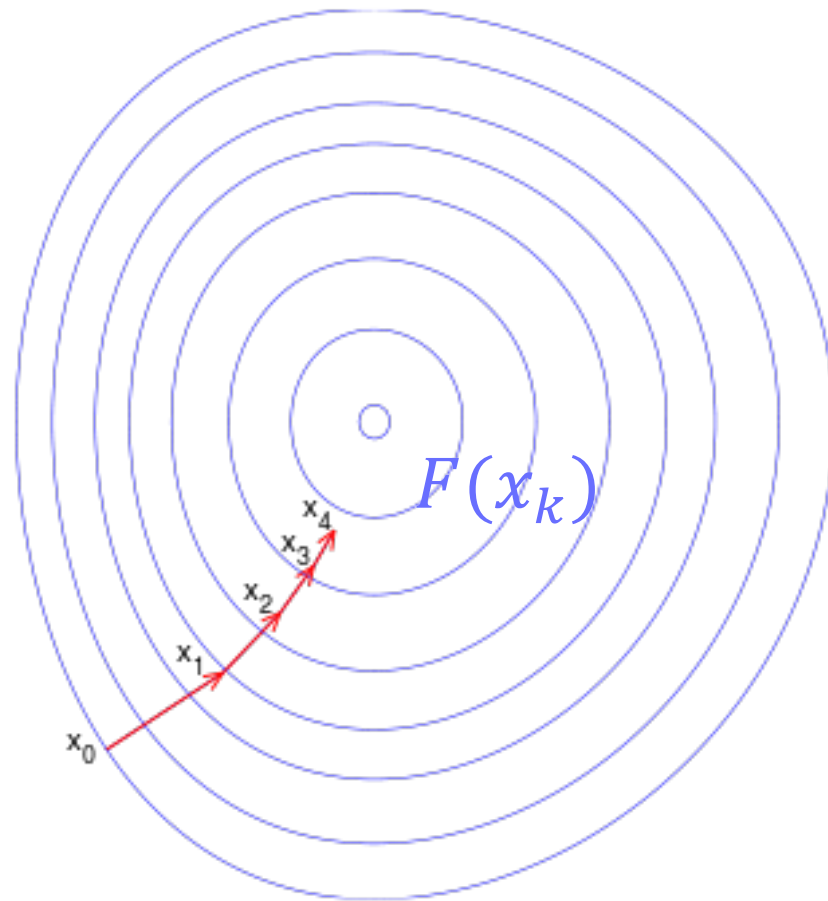
$$a(x) = \langle w, x \rangle + w_0$$

$$Q = \sum_{i=1}^l L(y_i, a(x_i)) \rightarrow \min_w$$

Обучение модели

Решение задачи оптимизации

Численное решение возможно с помощью
градиентного спуска (GD, Gradient Decent)



$$x_{k+1} = x_k - \gamma_k \nabla F(x_k)$$

Обучение модели

Стохастический градиент (SGD)

$$w_{k+1} = w_k - \gamma_k \sum_{i=1}^l y_i x_i L'(M_i)$$

$$w_{k+1} = w_k - \gamma_k y_i x_i L'(M_i)$$

x_i — случайный элемент обучающей выборки

Обучение модели

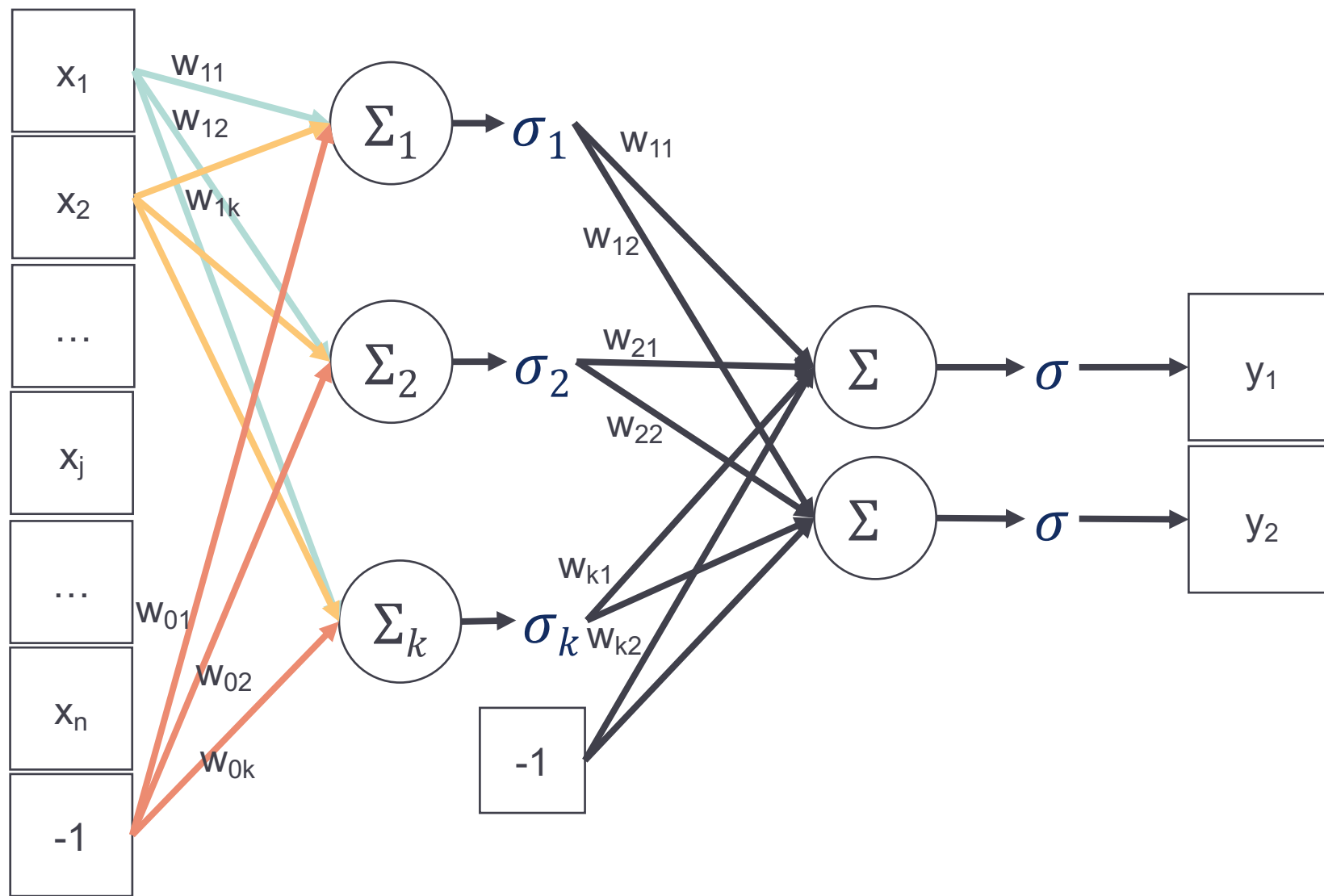
Применим SGD для нейронной сети

$$Q = \sum_{i=1}^l L(y_i, a(x_i)) \rightarrow \min_w$$

Проблема: в нейронной сети гораздо больше параметров w . Это приводит к проблемам при расчете градиента функции потерь $L'(w)$

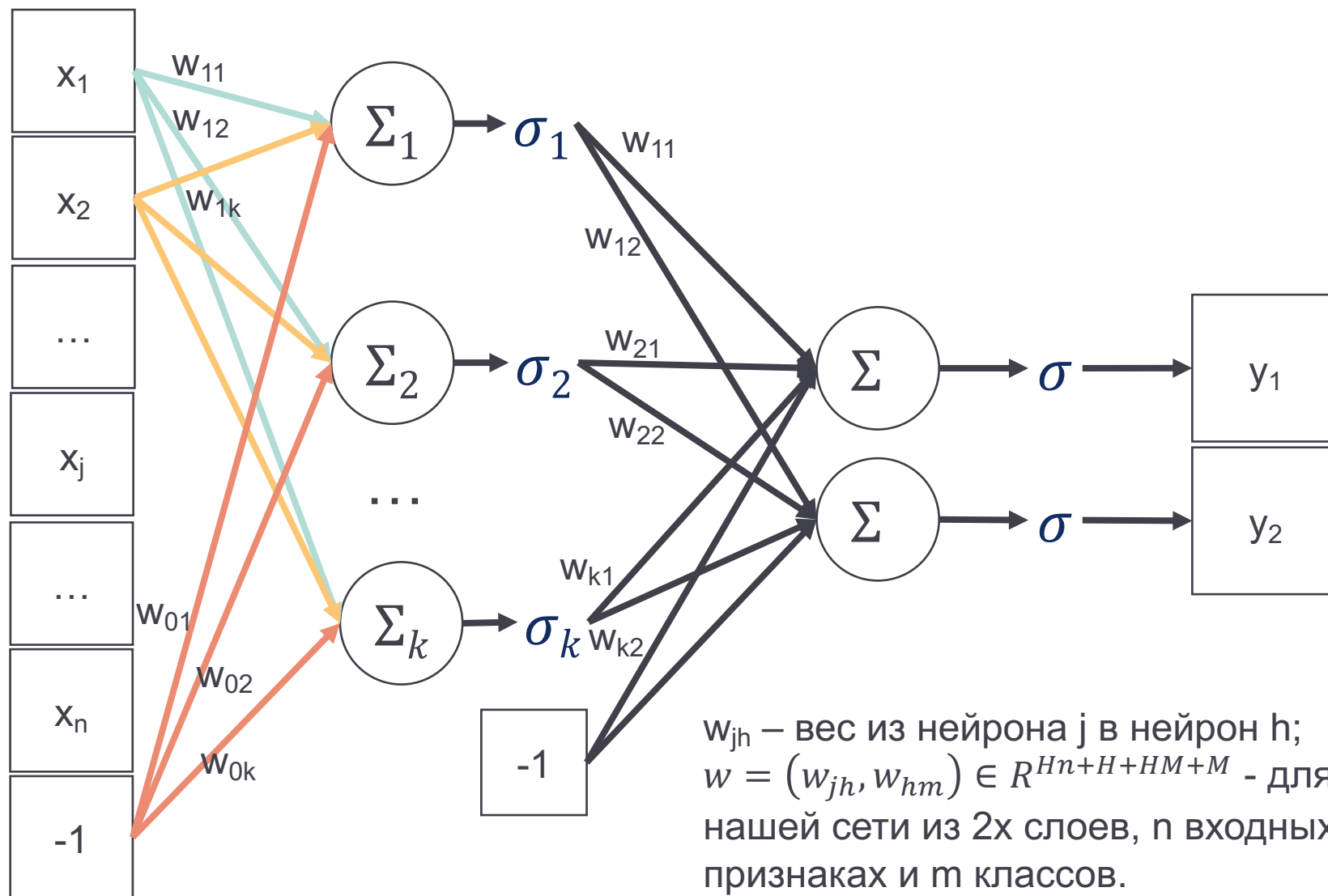
Обучение модели

Нейронная сеть



Обучение модели

Нейронная сеть



Back Propagation

Продолжим анализировать пример с 2х-слойной нейронной сетью.

Запишем прогноз нейрона m на последнем слое:

$$a^m(x_i) = \sigma_m(\sum_{h=0}^H w_{hm} u^h(x_i)), \text{ где}$$

$u^h(x_i) = \sigma_h(\sum_{j=0}^J w_{jh} f_j(x_i))$ - нейроны с промежуточного слоя.

Back Propagation

Рассмотрим квадратичную функцию потерь

$$L_i(w) = \frac{1}{2} \sum_{m=1}^M (a^m(x_i) - y_i^m)^2$$

На каждом объекте мы получаем прогноз в виде вектора размера M . Соответственно, потери мы считаем также по векторам.

Back Propagation

Чтобы оптимизировать веса градиентными методами, нужно научиться считать градиент потерь по всем весам.

Для этого далее потребуется оценить $\frac{\partial L_i(w)}{\partial a^m}$ и $\frac{\partial L_i(w)}{\partial u^h}$, давайте с них начнем:

$$\frac{\partial L_i(w)}{\partial a^m} = a^m(x_i) - y_i^m, \text{ обозначим } \varepsilon_i^m$$

$$\frac{\partial L_i(w)}{\partial u^h} = \sum_{m=1}^M \frac{\partial L_i(w)}{\partial a^m} \sigma'_m w_{hm} = \sum_{m=1}^M \varepsilon_i^m \sigma'_m w_{hm}, \text{ обозначим } \varepsilon_i^h$$

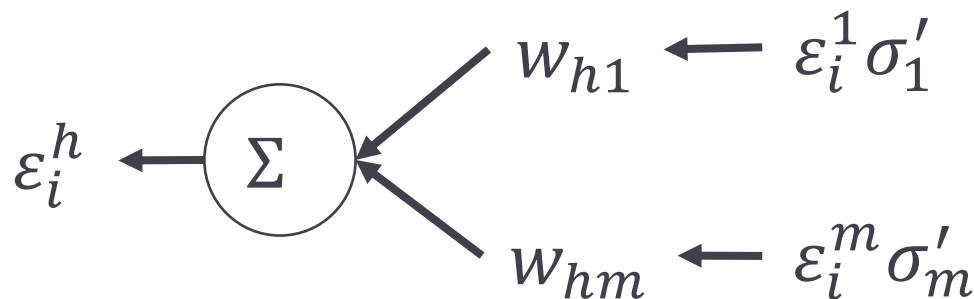
Назовём ε_i^h ошибкой на скрытом слое h

Back Propagation

$$\frac{\partial L_i(w)}{\partial u^h} = \sum_{m=1}^M \varepsilon_i^m \sigma'_m w_{hm} = \varepsilon_i^h$$

Получается, что ε_i^h вычисляется по ε_i^m если запустить сеть задом наперёд (backward pass)

Рассмотрим нейрон со скрытого слоя h:



Получилось, что мы пропустили ошибку обратным ходом через сеть. Отсюда back propagation.

Back Propagation

А получив частные производные L_i по a^m и по u^h легко найти производные и по весам.

На скрытом слое (по всем m и h):

$$\frac{\partial L_i(w)}{\partial w_{hm}} = \frac{\partial L_i(w)}{\partial a^m} \frac{\partial a^m}{\partial w_{hm}} = \varepsilon_i^m \sigma'_m u^h(x_i)$$

На входном слое (по всем h и j) :

$$\frac{\partial L_i(w)}{\partial w_{1h}} = \frac{\partial L_i(w)}{\partial u^h} \frac{\partial u^h}{\partial w_{jh}} = \varepsilon_i^h \sigma'_h f_j(x_i)$$

Обучение модели

Оптимизация модели методом Back Propagation

1. Инициализируем исходные веса
2. Выбираем объект или объекты для того, чтобы сделать шаг оптимизации.
3. На каждом шаге:
 - делаем прямой проход (forward pass), чтобы посчитать значения на всех скрытых слоях, прогноз и ошибки на выходном слое
 - делаем обратный проход (backward pass), чтобы рассчитать ошибки на скрытых слоях и оценить частные производные
 - зная производные, делаем шаг - обновляем веса
 - считаем потери
4. Повторяем шаг 3, пока потери не стабилизировались

Обучение модели

Оптимизация модели методом Back Propagation

1. Инициализируем исходные веса w
2. Выбираем объект x_i или объекты для того, чтобы сделать шаг оптимизации.
3. На каждом шаге:
 - делаем прямой проход (forward pass), чтобы посчитать значения на всех скрытых слоях $u^h(x_i) = \sigma_h(\sum_{j=0}^J w_{jh} f_j(x_i))$, прогноз $a^m(x_i) = \sigma_m(\sum_{h=0}^H w_{hm} u^h(x_i))$ и ошибки на выходном слое $\varepsilon_i^m = \frac{\partial L_i(w)}{\partial a^m}$
 - делаем обратный проход (backward pass), чтобы рассчитать ошибки на скрытых слоях и оценить частные производные $\varepsilon_i^h = \sum_{m=1}^M \varepsilon_i^m \sigma'_m w_{hm}$
 - зная производные, делаем шаг - обновляем веса $w_{jh} = w_{jh} - \eta \varepsilon_i^h \sigma'_h f_j(x_i)$ и $w_{hm} = w_{hm} - \eta \varepsilon_i^m \sigma'_m u^h(x_i)$
 - считаем потери $Q = (1 - \lambda)Q + \lambda L_i$
4. Повторяем шаг 3, пока потери не стабилизировались

Обучение модели

Back Propagation

Достоинства

- Быстрое вычисление градиента
- Потокное обучение, распределенное обучение
- Простое обобщение для большого числа слоёв, различных архитектур, функций потерь и функций активации

Ограничения

- Скорость сходимости
- Локальные экстремумы
- Переобучение
- Паралич сети (горизонтальные асимптоты функций активации)

Подробнее про Back Propagation

- Реализация на Python:
<https://machinelearningmastery.com/implement-backpropagation-algorithm-scratch-python/>
- Математическое изложение, см. стр. 13 в лекциях К.В. Воронцова по нейронным сетям:
<http://www.ccas.ru/voron/download/NeuralNets.pdf>
- Популярное и подробное изложение (english), см. стр. 200-217 в Deep Learning Book:
<https://www.deeplearningbook.org/contents/mlp.html>

Оптимизация модели

Оптимизация сходимости

Идея: можем ли мы улучшить способ обновления весов методом градиентного спуска, чтобы ускорить сходимость?

Оптимизация модели

Оптимизация сходимости

Метод накопления импульса:

- Momentum – экспоненциальное скользящее среднее градиента по последним итерациям
- NAG (Nesterov's accelerated gradient) – стохастический градиент с импульсом Нестерова

Адаптивный шаг:

- RMSProp (running mean square) – адаптация скорости изменения весов скользящим средними по последним итерациям
- AdaDelta (adaptive learning rate) – двойная нормировка приращений весов

А также другие методы и их комбинации!

Оптимизация модели

Пактеная нормализация данных (batch normalization)

$B = \{x_i\}$ – mini-batch или пакет данных

Оптимизация градиентов $L_i(w)$ по пакету ускоряет сходимость

Функции активации

Такие функции активации как сигмоид или гиперболический тангенс могут привести к затуханию градиентов

Функция положительной срезки (rectified linear unit):

- $\text{ReLU}(y) = \max\{0; y\}$
- $\text{PReLU}(y) = \max\{0; y\} + \alpha \min\{0; y\}$ (parametric ReLU)

Начальное приближение весов

Существуют различные эвристики для начального приближения весов:

- Выбор весов случайно из распределений с нулевым средним, например, нормального или равномерного
- Инициализация весами предобученной модели

Интересная идея для начального приближения - обучение нейронов как независимых линейных моделей:

- по подмножествам признаков
- по подмножествам объектов
- из разных начальных приближений

В этом случае обучение будет послойным. Далее модель обучается методом back propagation.

Регуляризация модели

Регуляризация модели

Идея: можем ли мы сделать нейронную сеть более устойчивой?

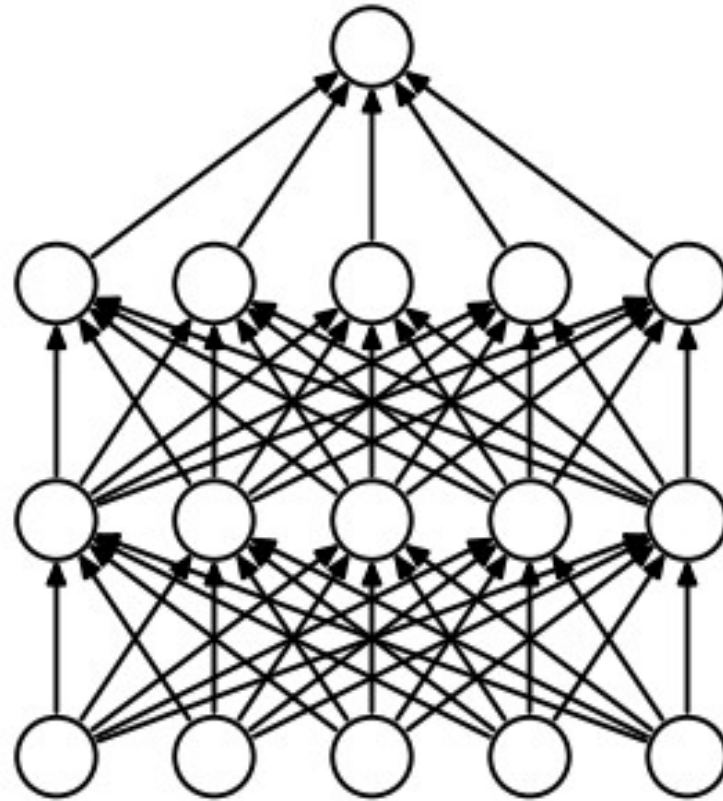
Регуляризация модели

Dropout

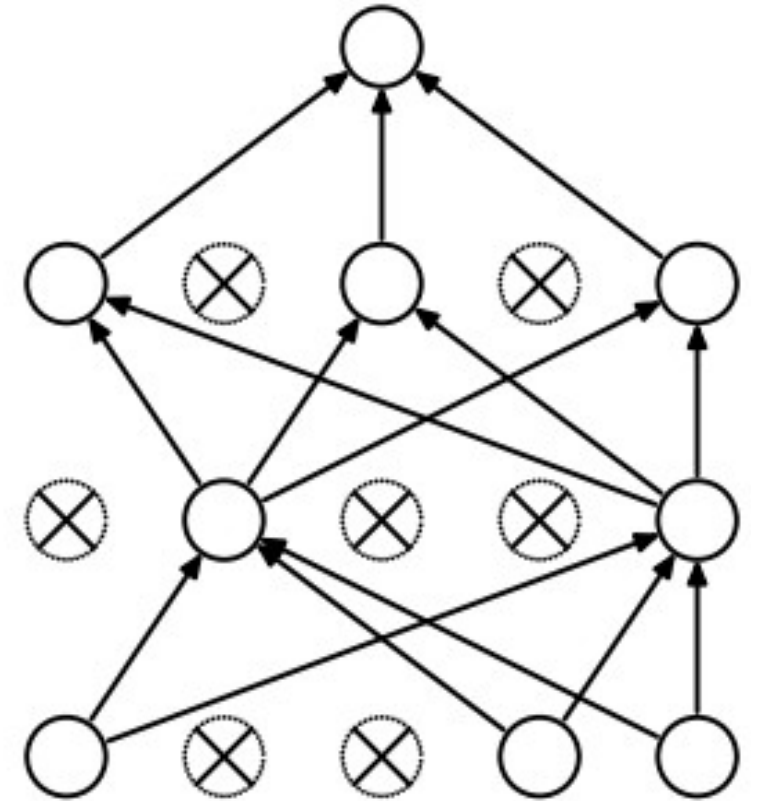
Идея: можем ли мы сделать нейронную сеть более устойчивой?

Более устойчивая сеть – такая сеть, которая при отключении части нейронов работает примерно также как исходная.

Dropout

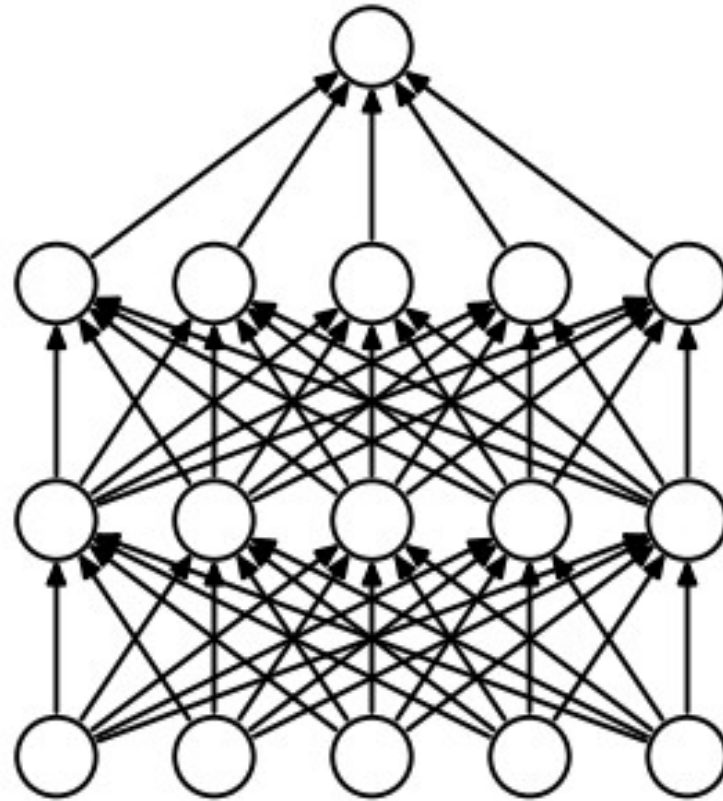


(a) Standard Neural Net

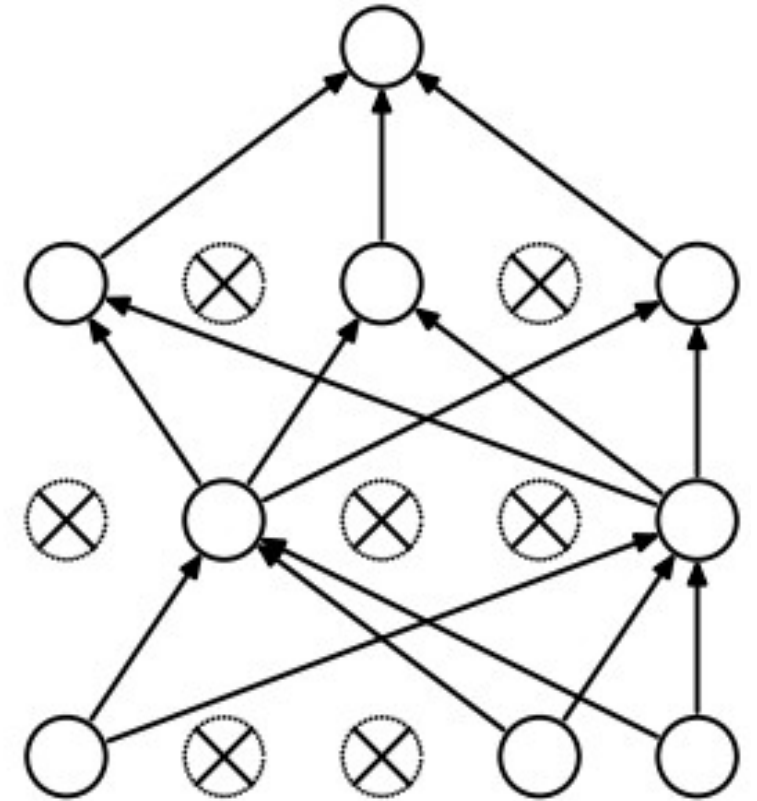


(b) After applying dropout.

Dropout



(a) Standard Neural Net



(b) After applying dropout.

С вероятностью p_l мы **отключаем h -ый** нейрон слоя l . Это означает, что на вход слоя $l+1$ из нейрона h предыдущего слоя ничего не придет (придет 0).

Регуляризация модели

Dropout

- Dropout позволяет сделать нейронную сеть сверх устойчивой за счет того, что разные подмножества нейросети работают примерно также, как вся нейросеть
- Этот эффект достигается за счет отключения нейронов во время обучения
- Получается, что мы как будто учим разные подмножества сети приближать целевую функцию независимо, вместо того, чтобы они компенсировали ошибки друг друга

Dropout + l_2 регуляризация

- Помните, мы делали l_2 регуляризацию для линейных моделей? А именно минимизировали не только ошибки, но и квадраты весов

$$V(w) = \|w\|_{l_2}^2 = \sum_{n=1}^d w_n^2$$

Dropout + l_2 регуляризация

- Помните, мы делали l_2 регуляризацию для линейных моделей? А именно минимизировали не только ошибки, но и квадраты весов

$$V(w) = \|w\|_{l_2}^2 = \sum_{n=1}^d w_n^2$$

- Этот подход можно комбинировать с Dropout. Изменится формула для расчета градиента потерь при обучении нейросети, там появится дополнительное слагаемое от регуляризации.

Популярные приложения

Популярные приложения

Популярные приложения

- Computer vision: классификация изображений, детекция объектов, генерация тестовых описаний изображений, обработка и стилизация и пр.
- Natural Language Processing: распознавание речи, оценка тональности текста, поддержание диалога и пр.
- Timeseries analysis: прогнозирование спроса, прогнозирование популяции животных и пр.

А также многое другое: обработка звука; обработка видео; генерация изображений, музыки, текстов и пр.

Популярные приложения

Популярные приложения

- Computer vision
- Natural Language Processing
- Timeseries analysis
- и др.

Для работы с каждым из типов задач используются специализированные архитектуры сетей и оптимизации процесса обучения.

Популярные приложения

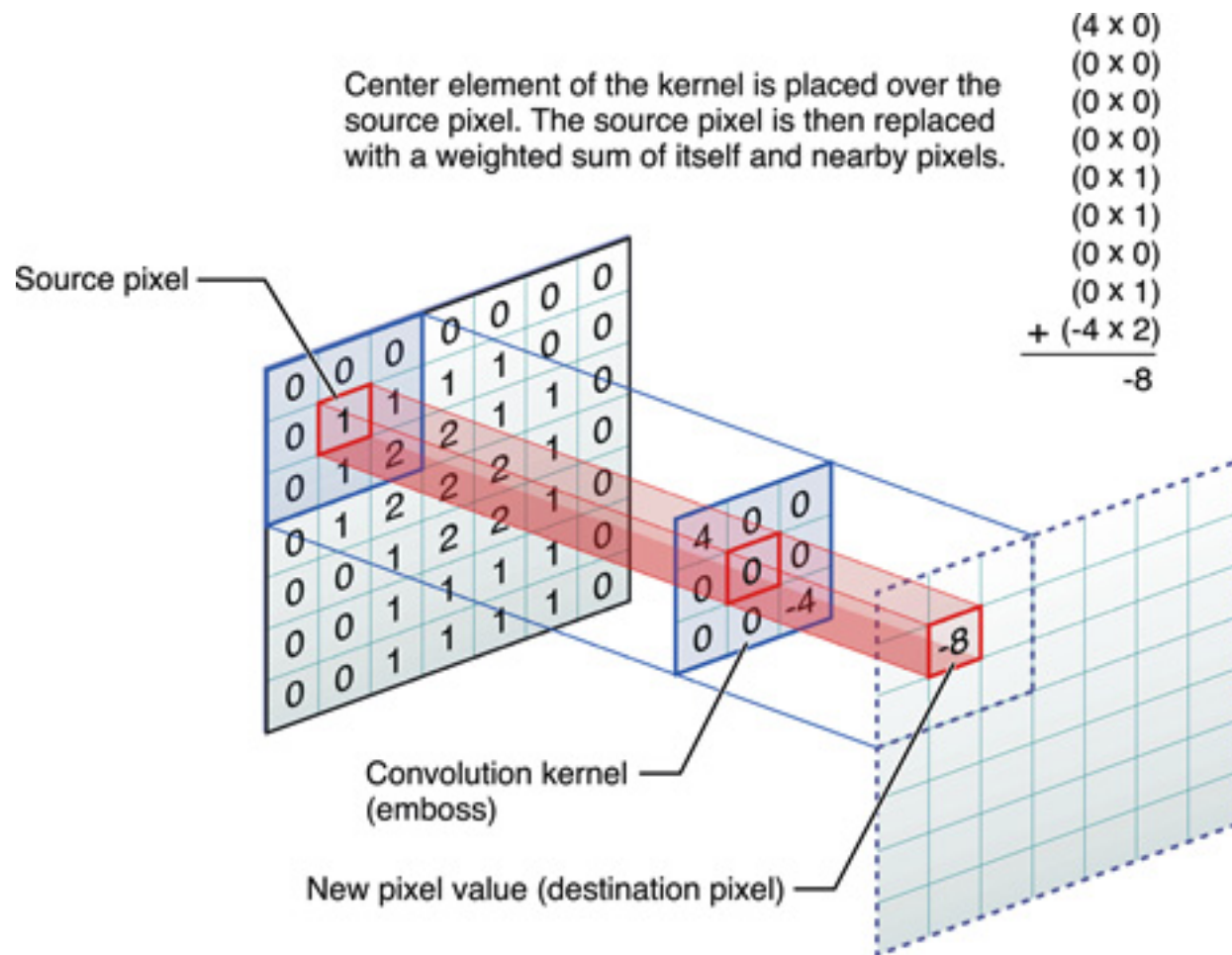
Для работы с каждым из типов задач используются специализированные архитектуры сетей и оптимизации процесса обучения.

Мы коротко обсудим:

- Сверточные сети
- Автоенкодеры
- Рекуррентные сети

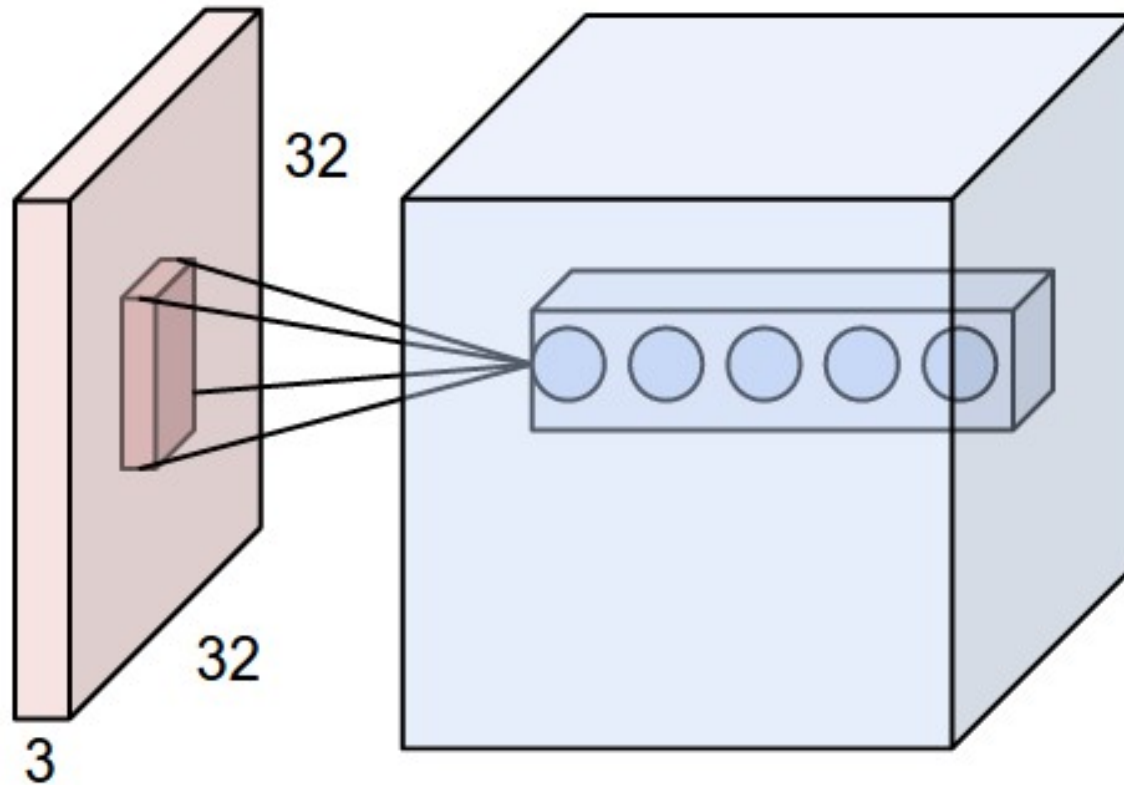
Convolution (свёртка)

Свёрточные сети



Convolution (свёртка)

Свёрточные
сети



Свёрточные сети

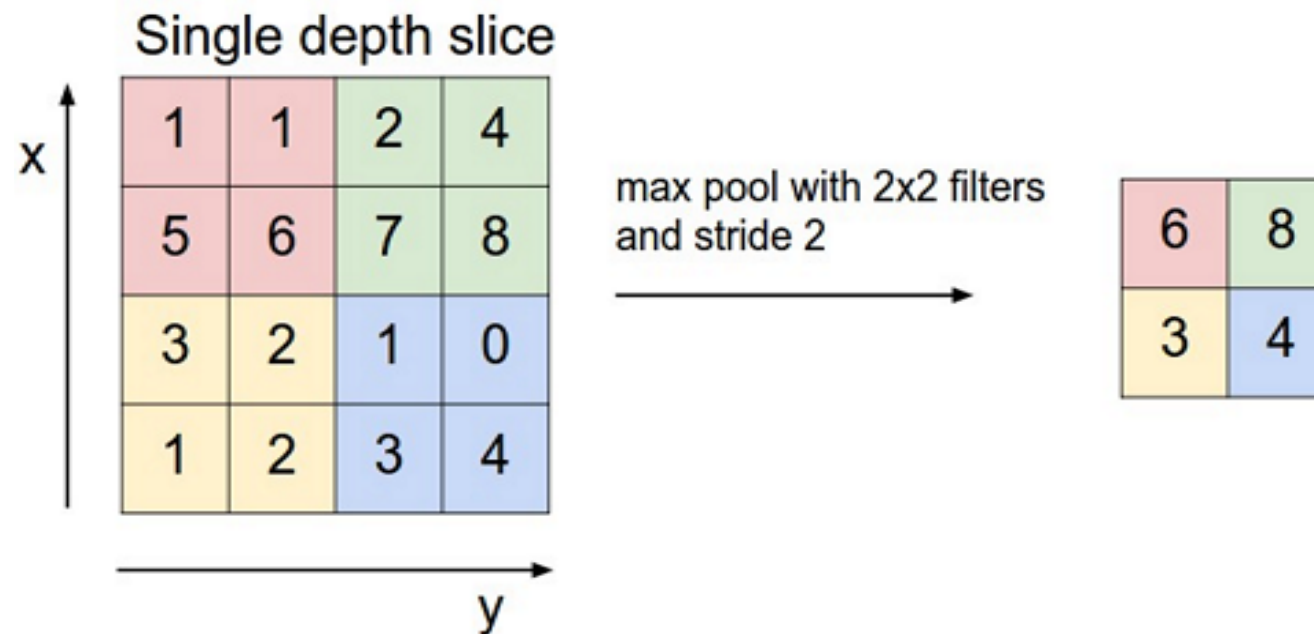
Pooling

Проблемы использования большего количества сверток:

- Размер изображения после свертки такой же или почти такой же, как был изначально
- Размерность не уменьшается, а еще фильтров будет много – значит только увеличивается
- Получается очень много параметров – есть риск не получить никакой обобщающей способности у алгоритма

Свёрточные сети

Pooling

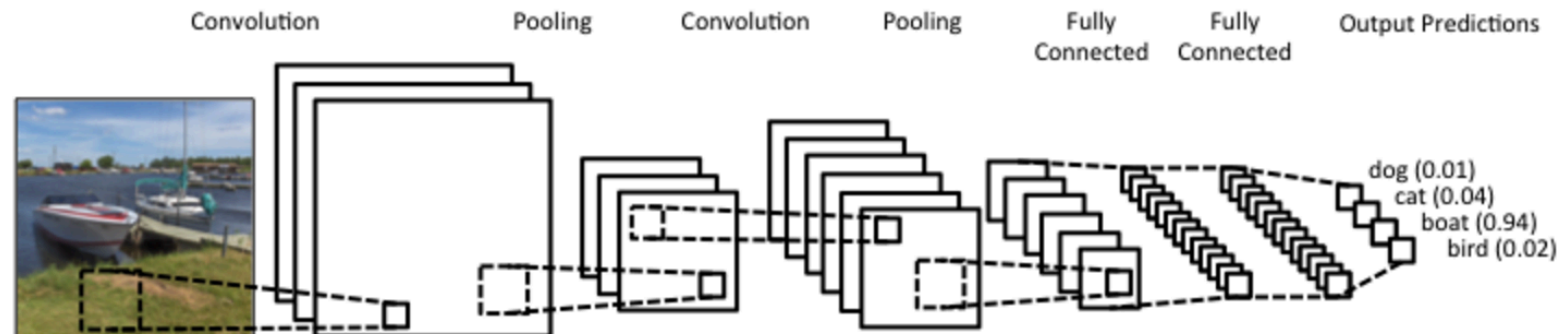


На практике используют

- Average Pooling
- Max Pooling

Сверточная нейронная сеть

Свёрточные сети



Пример: сеть GoogLeNet

Свёрточные сети

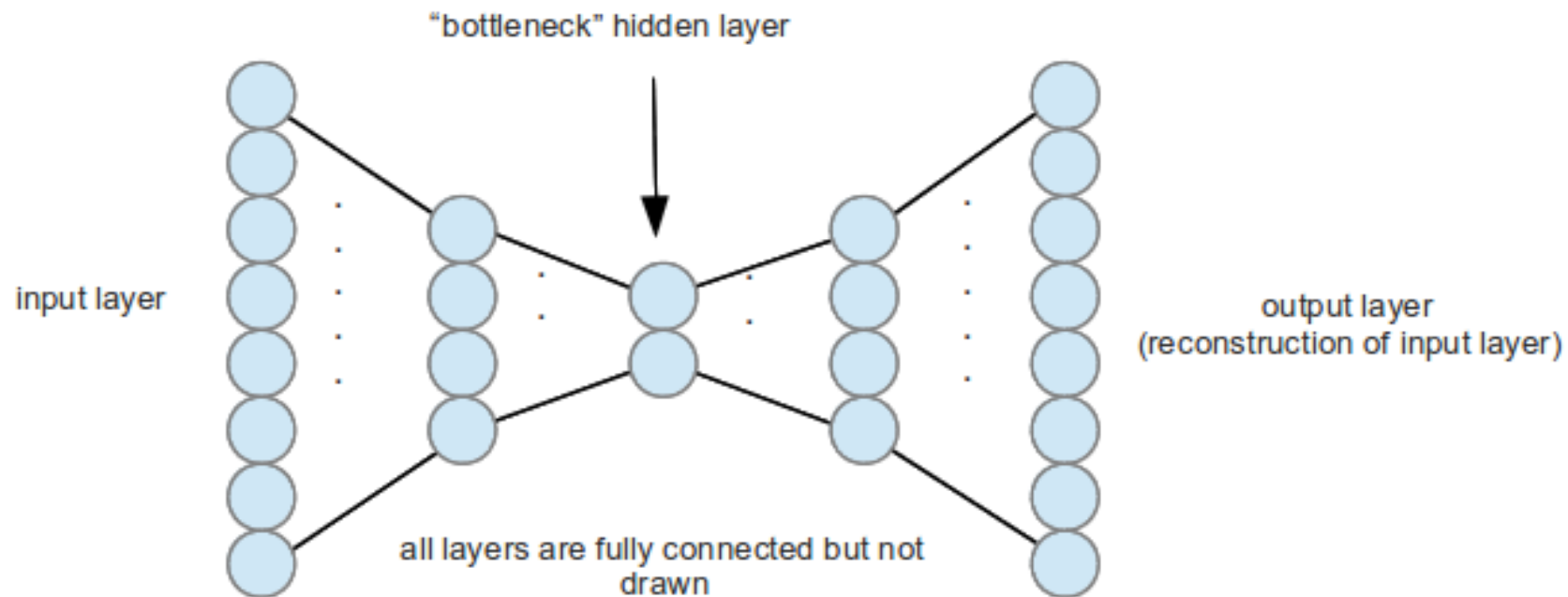


Convolution
Pooling
Softmax
Other

АВТОЭНКодеры

АВТОЭНКодеры

1. Построим первые слои нейросети так, чтобы в них было все меньше и меньше нейронов
2. Далее зеркально отразим слои

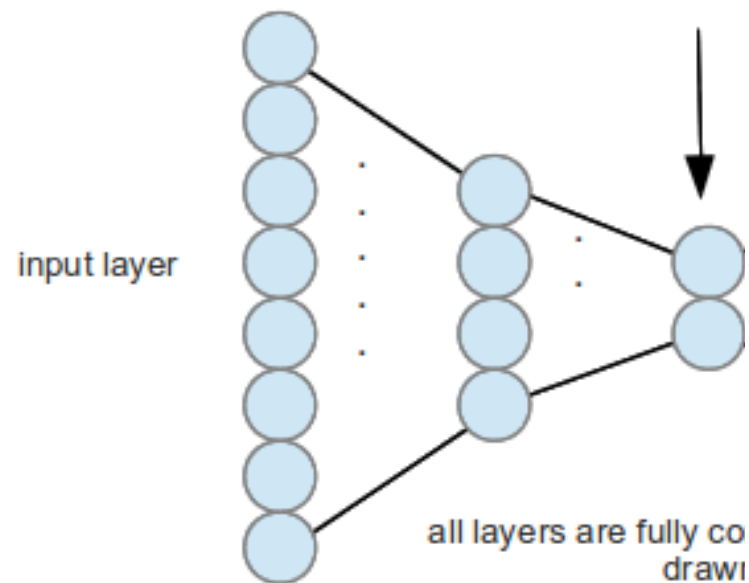


3. Обучим нейросеть восстанавливать на выходе то же, что было на входе.

Автоэнкодеры

Автоэнкодеры

4. «Срежем» всё после самой «узкой» части и будем использовать сеть для снижения размерности



АВТОЭНКОДЕРЫ

Виды автоенкодеров:

- Sparse autoencoders
- Denoising autoencoders
- Convolutional autoencoders

Автоэнкодеры

Автоэнкодеры

Как использовать полученные признаки?

- Обучать другие алгоритмы на них
- Обучать модель на них и исходных признаках
- Использовать автоэнкодеры для предобучения нейросетей, чтобы потом обучение сходилось к лучшему решению

Рекуррентные нейронные сети

При традиционном подходе к обучению нейронных сетей предполагается, что все входы и выходы независимы.

Данное предположение работает не всегда:

- NLP: мы хотим предсказать следующее слово предложения по k предыдущим
- Timeseries: мы хотим продолжить временной ряд

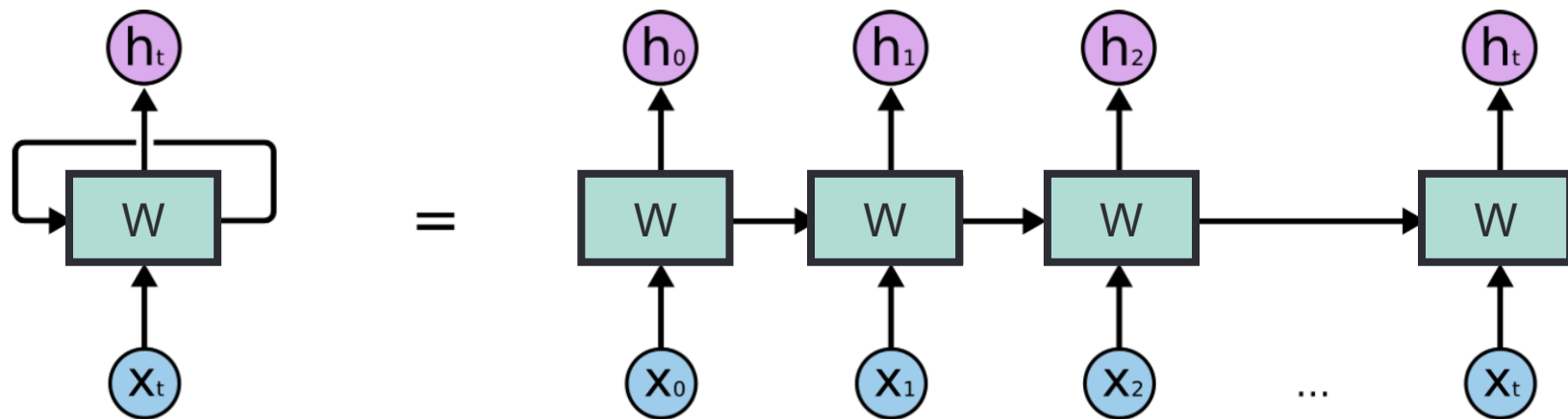
Существует много задач, где важно учитывать предыдущие входы. В терминах обучения сети: учитывать предыдущие состояния сети для вычисления текущего состояния.

Рекуррентные нейронные сети

Рекуррентные нейронные сети

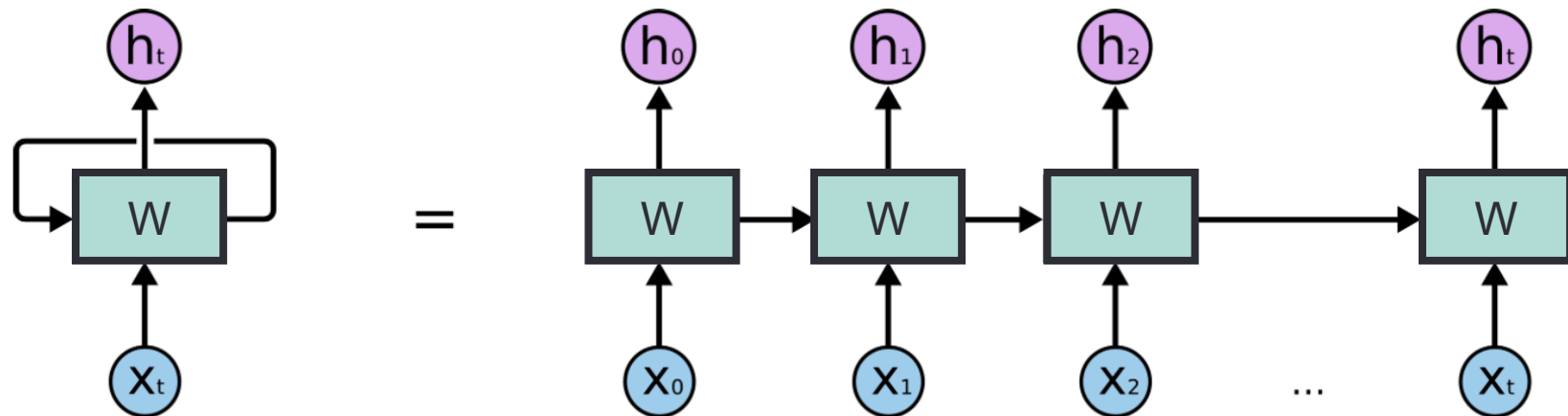
- Нейронная сеть преобразует слой t в слой $t + 1$ по правилу $h_{t+1} = \sigma_t(W_t h_t + b_t)$
- Получается, что нейросеть - это суперпозиция преобразований вида $a(x) = \sigma_t(W_t(\sigma_{t-1}(W_{t-1} \dots + b_{t-1})) + b_t)$

Идея: на каждом шаге обучения t будем обновлять посчитанные на предыдущем шаге веса W_t и b_t скрытого нейронного слоя



Рекуррентные нейронные сети

Идея: на каждом шаге обучения t будем обновлять посчитанные на предыдущем шаге веса W_t и b_t скрытого нейронного слоя

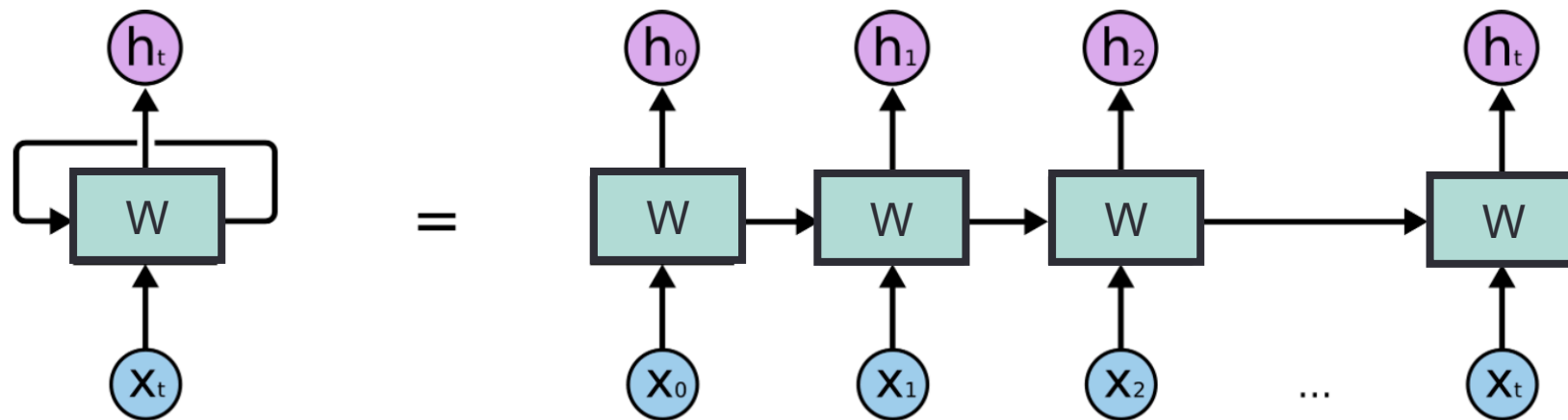


Тогда сеть будет задана следующим образом:

$$h_{t+1} = \sigma_t(W_{hh}h_t + W_{hx}x_t + b_h), \text{ где}$$

- Матрица W_{hh} определяет вклад вектора h с предыдущего шага в новый;
- Матрица W_{hx} определяет вклад нового объекта x .

Рекуррентные нейронные сети



Тогда сеть будет задана следующим образом:

$$h_{t+1} = \sigma_t(W_{hh}h_t + W_{hx}x_t + b_h), \text{ где}$$

- Матрица W_{hh} определяет вклад вектора h с предыдущего шага в новый;
- Матрица W_{hx} определяет вклад нового объекта x .

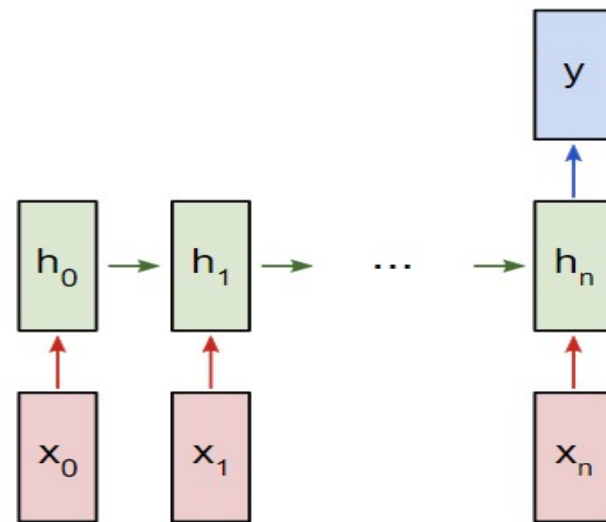
Осталось заметить, что на последнем слое нам нужно получить вектор ответов y , поэтому нужно добавить слой для вычисления выходного значения: $a(x) = \sigma_y(W_{hy}h_n + b_y)$

Рекуррентные нейронные сети

Рекуррентные нейронные сети

Таким образом простая рекуррентная сеть задается:

$$h_{t+1} = \sigma_t(W_{hh}h_t + W_{hx}x_t + b_h) \text{ и } a(x) = \sigma_y(W_{hy}h_n + b_y)$$



Здесь:

- Матрица W_{hh} определяет вклад вектора h с предыдущего шага в новый;
- Матрица W_{hx} определяет вклад нового объекта x ;
- Матрица W_{hy} и b_y позволяют вычислить выходное значение.

Рекуррентные нейронные сети

Простые нейронные сети можно усложнить для получения решений с лучшими свойствами (в частности более длительной памятью):

- LSTM (Long short-term memory)
- GRU (Gated Recurrent Unit)

Модель нейронной сети

Обучение нейронных сетей

- Модель персептрона и нейронная сеть известны давно
- Применение подхода back propagation позволило считать градиенты быстрее, но всё равно не достаточно быстро для обучения сетей с большим количеством скрытых слоев
- Современные вычислительные мощности и матричная модель вычислений позволили обучать нейросети с большим количеством слоев достаточно быстро
- Появились достаточно объемные наборы данных, для обучения моделей с большим количеством параметров: [Imagenet](#), [CIFAR-100](#)
- Сегодня появляются новые архитектуры нейросетей, позволяющие более качественно решать соответствующие прикладные задачи

Практические рекомендации

Как обучить нейронную сеть

Как обучить нейронную сеть

Библиотеки для обучения моделей:

- **Tensorflow**
- **Keras** (входит в Tensorflow, начиная с версии 2)
- **PyTorch**
- Theano
- Lasagne
- Nolearn
- Caffe

Как обучить нейронную сеть

Как обучить нейронную сеть

Для обучения модели потребуется:

1. Установить одну из библиотек на выбор
2. Найти любой пример с нейросетью на Keras/PyTorch/Tensorflow и воспроизвести у себя
3. Пробовать использовать разные архитектуры нейросетей (полносвязные, сверточные, рекуррентные, автоэнкодеры)

Как обучить нейронную сеть

Как обучить нейронную сеть

Будет полезно подробнее изучить Deep Learning.

- Стенфордский курс по computer vision:
<http://cs231n.stanford.edu/>
- Курс по NLP от deeplearning.ai
<https://www.deeplearning.ai/natural-language-processing-specialization/>
- Обучение и деплой моделей:
<https://course.fullstackdeeplearning.com/>

Курсов много, ищите курс, который вам понравится.

Машинное обучение: нейронные сети и глубокое обучение (введение)

Спасибо!
Эмили Драль