

Практический анализ данных и машинное обучение: искусственные нейронные сети

Ульянкин Филипп

23 октября 2018 г.

Автокодировщики, трюки, используемые при обучении сетей

Agenda

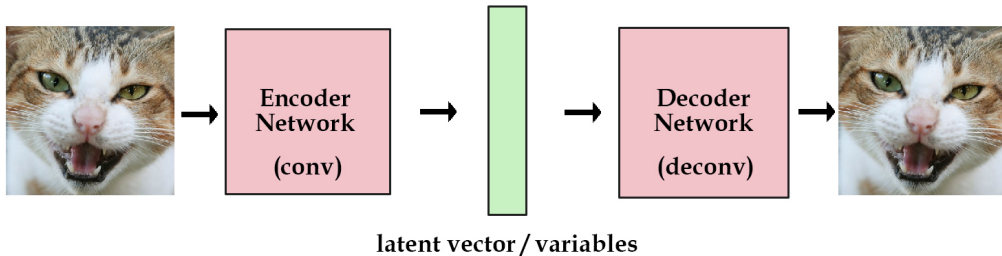
- Автокодировщики
- Эвристики, используемые при обучении нейронных сетей

Автокодировщики

Автокодировщики

- Понижение размерности --- задача обучения без учителя
- Давайте превратим её в обучение с учителем!

Автокодировщики



Собираем свой автокодировщик

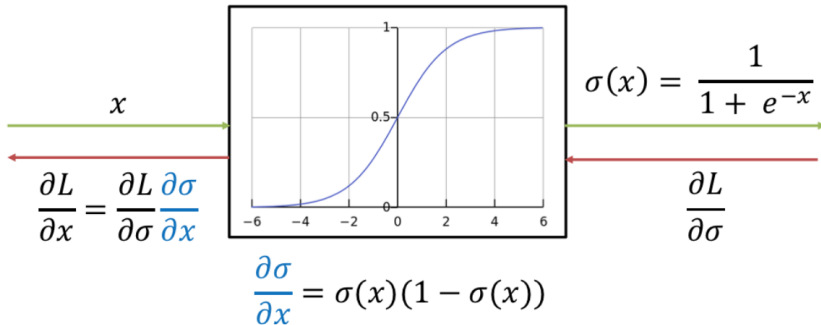
Эвристики для обучения сетей

Эвристики, используемые при обучении

- Применимы все те же эвристики, что и в градиентном спуске
- Инициализация весов, выбор шага, порядок предъявления объектов
- Кроме того появляются новые проблемы: выбор функции активации в каждом нейроне, выбор числа слоёв, числа нейронов, выбор значимых связей
- Это приводит к возникновению новых эвристик

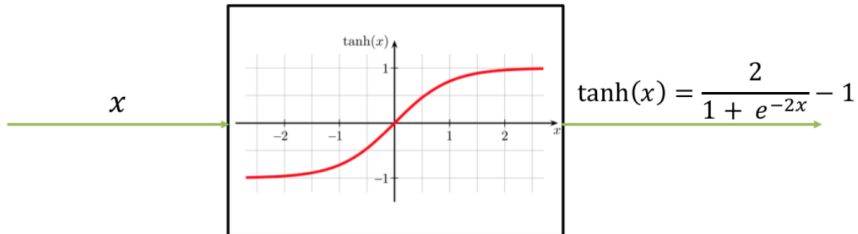
Функции активации

Sigmoid activation



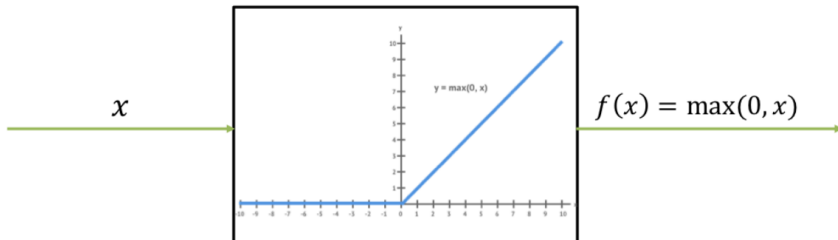
- В глубоких сетях способствует затуханию градиента (vanishing gradients)
- Не центрирована относительно нуля
- Вычислять e^x дорого

Tanh activation



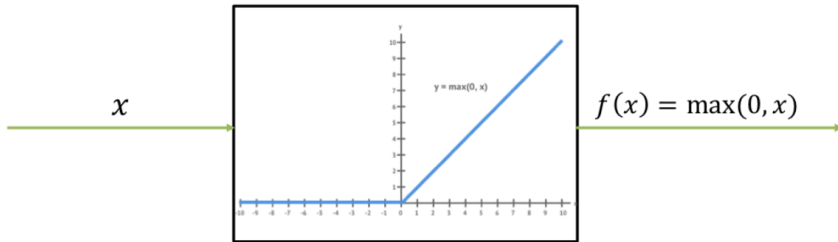
- Центрирован относительно нуля
- Всё ещё похож на сигмоиду

ReLU activation



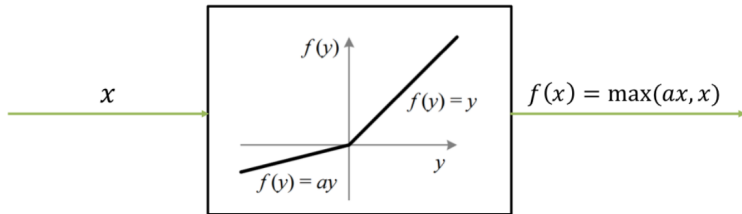
- Быстро вычисляется
- Градиент никогда не зануляется при $x > 0$
- Сходимость сеток ускоряется

ReLU activation



- Сетка может умереть, если активация занулится на всех нейронах
- Не центрирован относительно нуля

Leaky ReLU activation



- Как ReLU, но не умирает
- Важно, чтобы $a \neq 1$

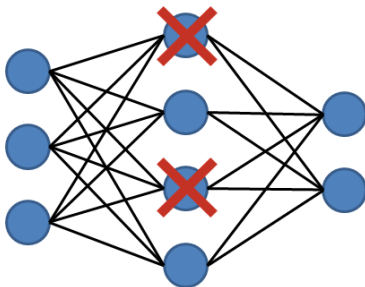
Что же выбрать

- На самом деле это неважно
- Важно собрать хорошую архитектуру и подобрать метод оптимизации
- Обычно в сетках работают либо с сигмоидом либо с ReLU, если остаётся время, то пробуют другие идеи, но обычно выигрыш в качестве от перебора в функциях активации довольно низкий

Дропаут

Dropout

- Придумали в 2014 году
- С вероятностью p отключаем нейрон
- Делает нейроны более устойчивыми к случайным возмущениям



Dropout

- На каждой итерации мы изменяем только часть параметров, обучение нейронов протекает чуть более независимо
- При тестировании все нейроны присутствуют в сетке, но их выходы домножаются на вероятность p
- На самом деле с байесовской точки зрения, вводя дропаут, мы обучаем ансамбль нейронных сетей

Регуляризация

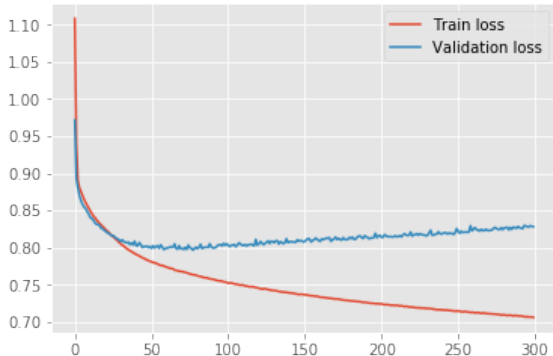
Регуляризация

- L_2 : приплюсовываем к функции потерь $\lambda \cdot \sum w_i^2$
- L_1 : приплюсовываем к функции потерь $\lambda \cdot \sum |w_i|$
- Можно регуляризовать не всю сетку, а отдельный нейрон или слой
- Не даёт нейрону сфокусироваться на слишком выделяющемся входе
- Очень похожа по своему действию на dropout

Early stopping

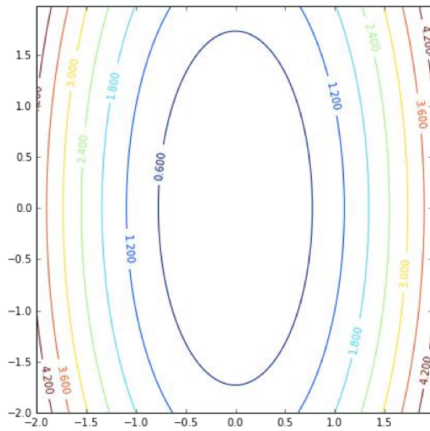
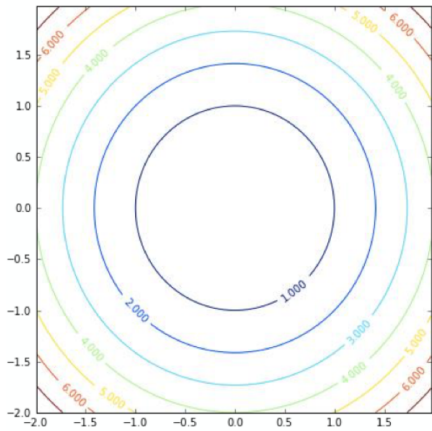
Early stopping

- Подбор числа шагов для обучения на основе проверки на валидации
- Для линейной модели с квадратичной (MSE) функцией потерь и SGD ранняя остановка эквивалентна L_2 регуляризации



Batchnorm

Стандартизация



Какая из ситуаций лучше для SGD?

Batch norm

- Придумали в 2015 году
- Давайте вместо X на входе использовать $\frac{X - \mu_X}{\sigma_X}$
- Давайте на каждом слое вместо $f(XW)$ использовать $\frac{f(XW) - \mu_f}{\sigma_f}$
- Обучение довольно сильно ускоряется

Инициализация модели

Инициализация весов

- Наши признаки X пришли к нам из какого-то распределения
- Выход слоя $f(XW)$ будет принадлежать другому распределению
- Если инициализировать веса неправильно, дисперсия распределения може от слоя к слою увеличиваться
- Эмпирически было выяснено, что это может портить сходимость для глубоких сетей

Инициализация весов

- Для симметричных функций с нулевым средним используйте инициализацию Ксавье `init="glorot_uniform"` или
- Для ReLU и им подобным инициализацию Хе `init="he_uniform"` или `init="he_normal"`
- Если вкратце, эти две инициализации корректируют параметры распределений в зависимости от входа и выхода слоя так, чтобы держать дисперсию в узде

Data augmentation

Data augmentation



Data augmentation

- Сдвиги
- Увеличение, уменьшение
- Повороты
- Искажение
- Затенение
- Смена стиля (красок)

Делает модель более устойчивой, полезна при маленьких выборках. На больших датасетах также улучшает результаты.

Другие приёмы

Предобучение

- Обучаем каждый нейрон на рандомной подвыборке, каждый нейрон впитает какие-то отдельные её особенности, после скрепляем все нейроны вместе и продолжаем обучение на всей выборке
- Обучаем на корпусе картинок автокодировщик, encoder благодаря этому учится выделять наиболее важные фичи, которые позволяют эффективно сжимать изображения. После срезаем decoder и на его месте достраиваем слои для решения нашей задачи, запускаем обычное дообучение.

Динамическое наращивание сети

- Обучение сети при заведомо недостаточном числе нейронов H
- После стабилизации функции потерь --- добавление нового нейрона и его инициализация путём обучения
 - либо по случайной подвыборке
 - либо по объектам с наибольшими значениями потерь
 - либо по случайному подмножеству входов
 - либо из различных случайных начальных приближений
- Снова итерации BackProp

Эмпирический опыт: Общее время обучения обычно лишь в 1.5 — 2 раза больше, чем если бы в сети сразу было итоговое число нейронов. Полезная информация, накопленная сетью не теряется при добавлении нейронов.

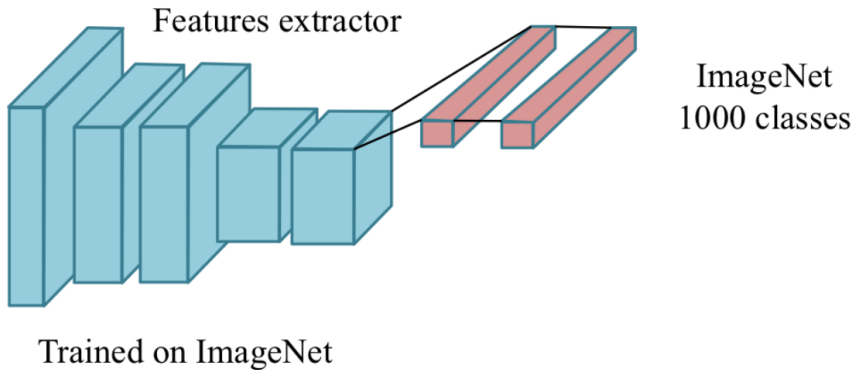
Прореживание сети

- Начать с большого количество нейронов и удалять незначимые по какому-нибудь критерию
- Пример: обнуляем вес, смотрим как сильно упала ошибка, сортируем все вязы по этому критерию, удаляем N наименее значимых
- После прореживания снова запускаем backprop
- Если качество модели сильно упала, вернуть последние удалённые связи

Transfer learning

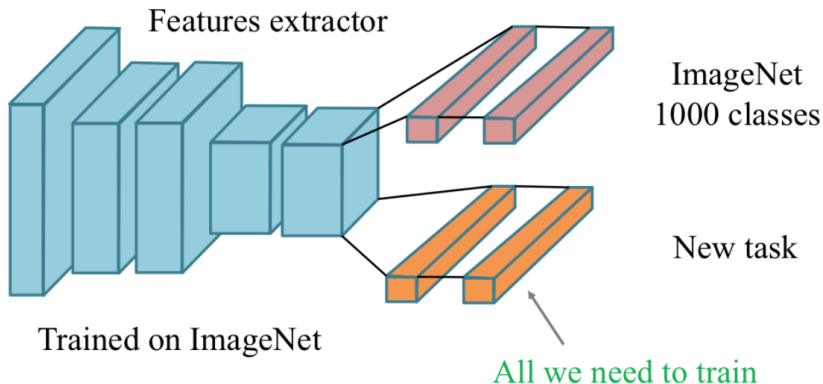
Transfer learning

- Глубокие сети извлекают из изображений сложные фичи, но для их обучения нужно много данных...



Transfer learning

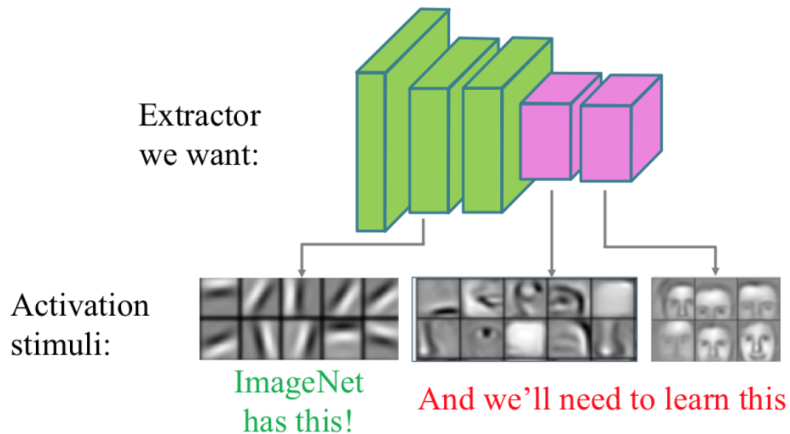
- Глубокие сети извлекают из изображений сложные фичи, но для их обучения нужно много данных...
- Давайте повторно использовать уже предобученную сеть!



Transfer learning

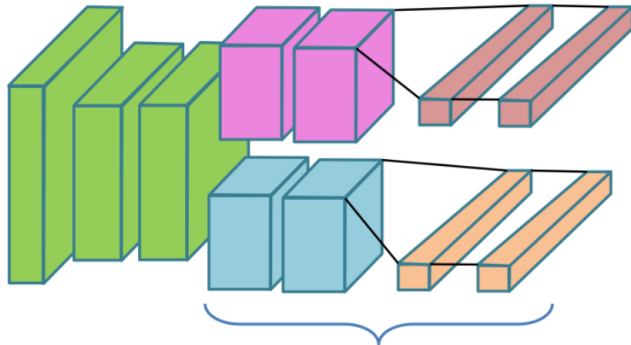
- Нужно меньше данных для обучения, так как нас интересуют лишь последние слои
- Это работает если наша задача похожа на ту, для которой обучалась используемая сетка
- Например, если мы хотим распознавать эмоции, в датасете для нашей сетки должны были быть человеческие лица

Transfer learning



Transfer learning

ImageNet features extractor



ImageNet
1000 classes

New task

All we need to train

Собираем свой transfer learning!