

Deep Learning. Улучшение нейронных сетей

Урок 2.5

Егор Конягин

1. Повторение
2. Проблема переобучения. Регуляризация
3. Методы регуляризации
4. Борьба с переобучением. Продолжение
5. Затухание/взрыв градиентов

Повторение

Задача оптимизации. Повторение

Основными проблемами, с которыми сталкиваются алгоритмы оптимизации при обучении нейронных сетей - это

- седловые точки;
- немасштабированные данные.

Для решения первой проблемы применяют алгоритмы с использованием скользящим средним градиентов предыдущих шагов (Momentum GD, AdaM), для решения второй: используют адаптивные алгоритмы (RMSProp, AdaM). Заметим, что AdaM является комбинацией адаптивных алгоритмов и алгоритмов со скользящим средним.

Важно! Для подсчета градиентов (не обновления параметров) используется SGD или mini-batch GD.

Проблема переобучения. Регуляризация

Проблема переобучения

Переобучение - это чрезмерное подстраивание под данные обучающей выборки, при котором ухудшается качество работы модели.



Рис. 1: Проблема переобучения

Переобучение или недообучение

Важно понимать, насколько хорошо данная модель обучена в текущий момент. Возможны три состояния:

1. оптимальное обучение;
2. недообучение (underfitting);
3. переобучение (overfitting).

Underfitting заключается в следующем: мы имеем похожее качество на train и на test - выборке, однако это качество нас не устраивает.

Overfitting характеризуется тем, что качество на train-выборке очень высокое (вплоть до 100%), в то время как на test-выборке оно может быть существенно ниже

Кривая обучения

Кривая обучения - это график зависимости функции потерь от номера итерации обучения.

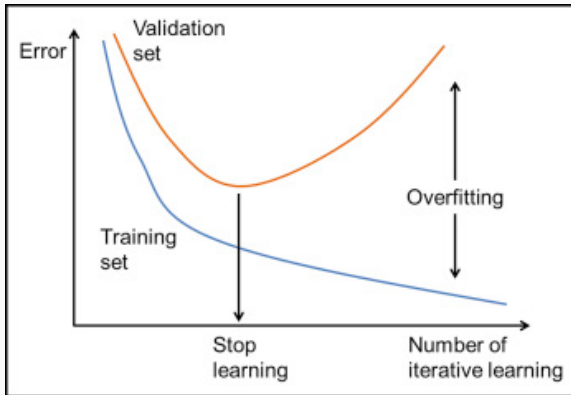


Рис. 2: Кривая обучения

Методы регуляризации

Регуляризация - это наложение неких ограничений на параметры w, b с тем, чтобы избежать переобучение. Сегодня мы рассмотрим следующие алгоритмы регуляризации:

1. L2 - регуляризация (ridge regularization);
2. L1 - регуляризация (Tikhonov, lasso regularization);
3. Dropout;

Weight decay

При проведенной нормализации входных данных (т.е. приведение их к виду $[-1, 1]$) одним из симптомов переобучения являются большие по абсолютному значению веса. Идея метода weight decay: увеличивать функцию потерь, если значения весов велики.

L2-регуляризация:

$$L(y, \hat{y})^{reg} = \sum \mathcal{L}(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2} \sum_l \|w^{[l]}\|_F^2 \quad (1)$$

L1-регуляризация (weight decay):

$$L(y, \hat{y})^{reg} = \sum \mathcal{L}(y^{(i)}, \hat{y}^{(i)}) + \lambda \sum_l \sum_{i,j} \|w_i^{[l]j}\|_F \quad (2)$$

Теперь градиент функции ошибки по параметрам изменится:

$$\frac{\partial L}{\partial w^{[l]}} + \frac{\lambda}{m} w^{[l]}. \quad (3)$$

Dropout

Dropout - это техника случайного отключения заданного количества нейронов с тем, чтобы увеличить стабильность работы нейросети.

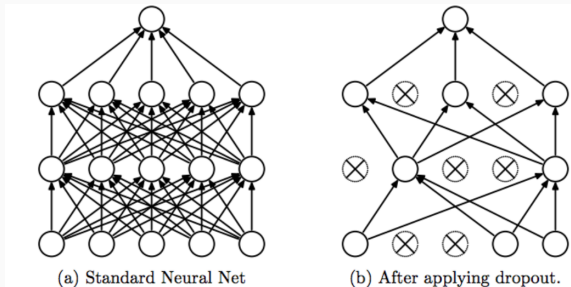


Рис. 3: Применение dropout. Источник: Stanford University

Dropout. Forward propagation

Рассмотрим, как же можно в формульном виде записать Dropout. Раньше forward propagation описывался так:

$$a^{[L]} = \sigma(W^{[L]}a^{[L-1]} + b^{[L]}) \quad (4)$$

Теперь надо выход каждого слоя поэлементно умножать на вектор, состоящий только из 1 и 0:

$$a^{[L]*} = a^{[L]} * D^{[L]}. \quad (5)$$

В итоге, в полученном выходе какие-то координаты будут равны нулю, а какие-то останутся без изменения.

Dropout характеризуется одним параметром - вероятностью.

Dropout. Backward propagation

Здесь также необходимо внести изменения: при подсчете $dA^{[L]}$ необходимо в конце этот вектор поэлементно умножить на маску $D^{[L]}$. Для этого надо не забыть сохранить ее значение, полученное в ходе forward propagation.

Помимо этого, вектор градиента надо разделить на параметр dropout-а, чтобы среднее значение выхода этого слоя оставалось неизменным.

ВАЖНО!!! Dropout применяется ТОЛЬКО на стадии обучения модели. Если модель уже обучена, dropout применять не нужно.

Проблема переобучения

Переобучение - это чрезмерное подстраивание модели под обучающие данные, при котором ухудшается качество работы модели.



Рис. 4: Проблема переобучения

Переобучение можно по-разному трактовать. Одна из возможных трактовок - это попытка нейросети "запомнить" правильные ответы, а не пытаться подстроиться под общую закономерность.

В результате, на тренировочных данных модель запоминает правильные ответы, и качество стремится к 100%. На тестовых данных модель, игнорируя реальные физические закономерности, выдает ответы, близкие к случайным, или просто с очень большой ошибкой.

Кривая обучения

Кривая обучения - это график зависимости функции потерь от номера итерации обучения.

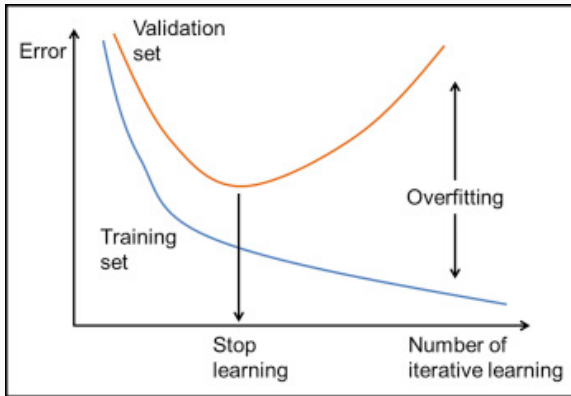


Рис. 5: Кривая обучения

Dropout

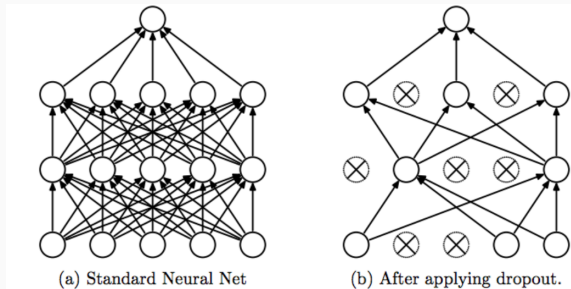


Рис. 6: Применение dropout. Источник: Stanford University

Борьба с переобучением.

Продолжение

Как мы ранее выяснили, борьба с переобучением осуществляется с помощью регуляризации, то есть в наложении дополнительных ограничений на функцию потерь.

Вопрос Может ли увеличение объема датасета побороться с переобучением?

А с недообучением?

Увеличение объема датасета иногда помогает бороться с переобучением.

Проблема не решается, если обучение алгоритма проходит на хороших данных, а работать нейросети приходится на худших данных (и наоборот). Это вызывает "перекосы" в распределении данных.

Затухание/взрыв градиентов

Рассмотрим глубокую нейросеть ($L \gg 1$):

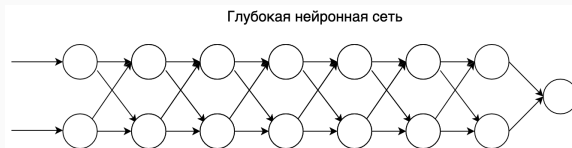


Рис. 7: Глубокая нейронная сеть

Положим, что все веса при инициализации являются положительными и меньше единицы. Выберем функцию активацию такую, что она удовлетворяет следующему разложению вблизи малого значения аргумента: $\sigma(x) = x + o(x)$.

Затухание/взрыв градиентов

Если это так, то forward propagation можно представить следующим способом:

$$\hat{y} = w^{[L]} \cdot w^{[L-1]} \dots w^{[1]} \cdot x. \quad (6)$$

По предположению веса инициализированы положительными малыми числами. Тогда

$$|y| \leq \|w\|_{max}^L \cdot |x|_{max} \quad (7)$$

Поскольку веса малы, то $\|w\| \rightarrow 0|_{L \rightarrow +\infty}$. Таким образом, у будет очень малым значением, обучение будет происходить крайне долго. Если веса, наоборот, будут больше единицы, то у будет экспоненциально возрастать как функция от L.

Затухание/взрыв градиентов

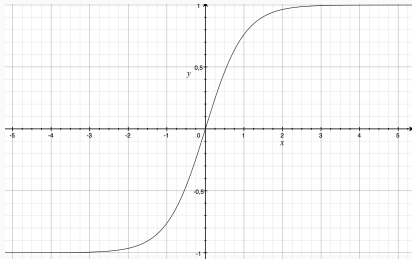


Рис. 8: График гиперболического тангенса

Как видно,

$$\tanh(x)|_{x \gg 1} \rightarrow 1; \quad \frac{d \tanh(x)}{dx}|_{x \gg 1} \rightarrow 0. \quad (8)$$

Таким образом, градиент такой функции активации в случае большого значения аргумента будет равен нулю. Обучение также будет происходить крайне долго.

Борьба с взрывом и затуханием градиентов. Инициализация весов

Инициализировать веса можно следующим образом (He initialization):

$$w^{[n_l]} \sim \mathcal{N}(0, 1) \cdot \sqrt{\frac{1}{n_{l-1}}} \quad - \text{sigmoid}; \quad (9)$$

$$w^{[n_l]} \sim \mathcal{N}(0, 1) \cdot \sqrt{\frac{2}{n_{l-1}}} \quad - \text{ReLU}. \quad (10)$$

Модификация

$$w^{[n_l]} \sim \text{Uniform}(-1, 1) \cdot \sqrt{\frac{2}{n_{l-1} + n_l}} \quad (11)$$

называется Xavier Initialization.

Batch normalization

Другой способ борьбы с взрывом и затуханием градиентов - это batch normalization (batch norm). Идея состоит в следующем: отнормировать данные на каждом слое:

$$\mu^{[l]} = \frac{1}{m} \sum_{i=1}^m z^{[l](i)} \quad (12)$$

$$\sigma^{[l]} = \frac{1}{m} \sum_{i=1}^m (z^{[l](i)} - \mu^{[l]}) \quad (13)$$

$$z_{norm}^{[l]} = \frac{z^{[l]} - \mu^{[l]}}{\sqrt{\sigma^2 + \varepsilon}} \quad (14)$$

Можно еще сильнее модифицировать z:

$$\hat{z}^{[l]} = \gamma^{[l]} z_{norm}^{[l]} + \beta^{[l]} \quad (15)$$

Сегодня мы

- поговорили про переобучение;
- познакомились с проблемой затухания градиентов;
- научились бороться с этой проблемой двумя методами
 1. правильной инициализацией весов;
 2. с помощью метода batch norm;