

Deep Learning. Оптимизация и настройка нейронных сетей - II

Урок 5

Егор Конягин

МФТИ & АО "ЦОСиВТ"

1. Повторение
2. Борьба с переобучением. Продолжение
3. Затухание/взрыв градиентов
4. Метод главных компонент (principal component analysis)

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

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



Рис. 1: Проблема переобучения

Кривая обучения

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

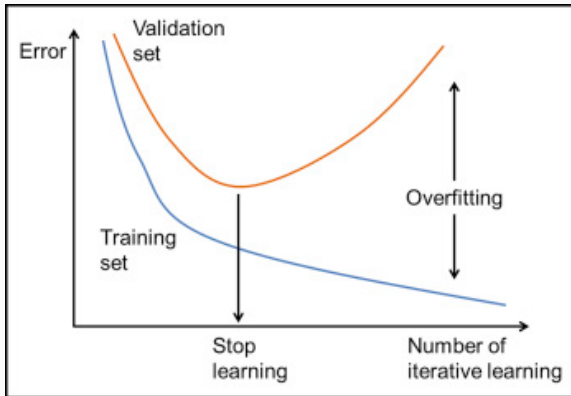


Рис. 2: Кривая обучения

Dropout

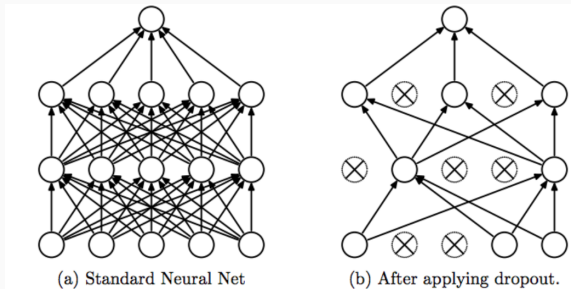


Рис. 3: Применение dropout. Источник: Stanford University

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 действительно помогает бороться с переобучением, поскольку с увеличением его объема распределение данных в обоих датасетах становятся более похожими.

Проблема иногда случается, что обучение алгоритма проходит на хороших данных, а работать нейросети приходится на худших данных. Это вызывает проблемы в качестве, которые сложно непосредственно отнести к переобучению.

Затухание/взрыв градиентов

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

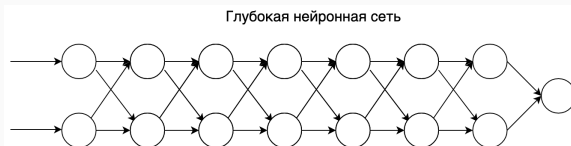


Рис. 4: Глубокая нейронная сеть

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:

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

По предположению веса инициализированы положительными малыми числами. Тогда

$$|y| \leq ||w||_{max}^L \cdot |x|_{max} \quad (2)$$

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

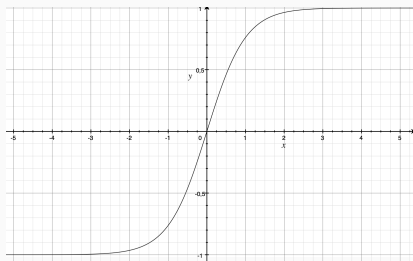


Рис. 5: График гиперболического тангенса

Как видно,

$$\tanh(x)|_{x \gg 1} \rightarrow 1; \quad \frac{d \tanh(x)}{dx}|_{x \gg 1} \rightarrow 0. \quad (3)$$

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):

$$w^{[n_l]} \sim \mathcal{N}(0, 1) \cdot \sqrt{\frac{1}{n_{l-1}}} \quad - \text{sigmoid}; \quad (4)$$

$$w^{[n_l]} \sim \mathcal{N}(0, 1) \cdot \sqrt{\frac{2}{n_{l-1}}} \quad - \text{ReLU}. \quad (5)$$

Модификация

$$w^{[n_l]} \sim \mathcal{N}(0, 1) \cdot \sqrt{\frac{2}{n_{l-1} + n_l}} \quad (6)$$

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:

$$\mu^{[l]} = \frac{1}{m} \sum_{i=1}^m z^{[l](i)} \quad (7)$$

$$\sigma^{[l]} = \frac{1}{m} \sum_{i=1}^m (z^{[l](i)} - \mu^{[l]}) \quad (8)$$

$$z_{norm}^{[l]} = \frac{z^{[l]} - \mu^{[l]}}{\sqrt{\sigma^2 + \varepsilon}}. \quad (9)$$

Можно еще сильнее модифицировать z :

$$\hat{z}^{[l]} = \gamma^{[l]} z_{norm}^{[l]} + \beta^{[l]} \quad (10)$$

Метод главных компонент (principal component analysis)

Метод главных компонент (principal component analysis)

Метод главных компонент - это метод понижения размерности признакового пространства. РСА непосредственно не относится к глубокому обучению, однако его можно применять и в задачах глубокого обучения в рамках этапа предобработки данных.

Сингулярное разложение:

$$\mathbf{M} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*, \quad (11)$$

где \mathbf{U} —ортогональная (унитарная) матрица, $\mathbf{\Sigma}$ —диагональная матрица, причем кол-во ненулевых элементов на диагонали равно рангу матрицы \mathbf{M} , а \mathbf{V} —тоже ортогональная (унитарная) матрица (звездочка означает эрмитово сопряжение).

1. Нормализация:

$$X^* = \frac{X - \mu}{\sigma} \quad (12)$$

2. поиск первой главной компоненты (поиск направления, вдоль которого дисперсия максимальна):

$$\frac{1}{m} \sum_{i=1}^m (u^T X^{(i)})^2 = \frac{1}{m} \sum_{i=1}^m (u^T X^{(i)T} X^{(i)} u) = u^T (X^T X) u; \quad (13)$$

3. Поиск остальных главных компонент:

$$\hat{X}_k = X - \sum_{s=1}^{k-1} X u_{(s)} u_{(s)}^T; \quad (14)$$

$$u_k = \underset{u}{\operatorname{argmax}} (u^T \hat{X}_k^T \hat{X}_k u). \quad (15)$$

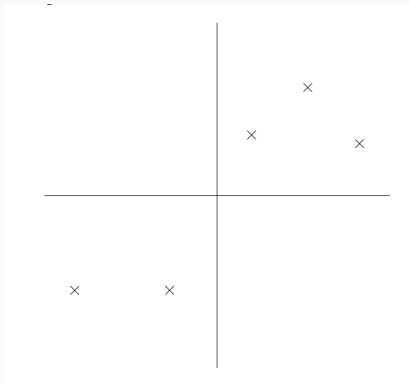
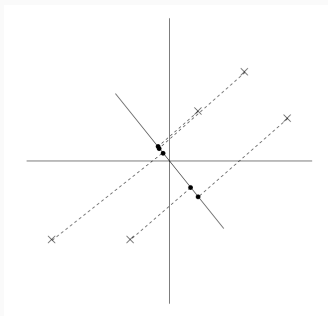
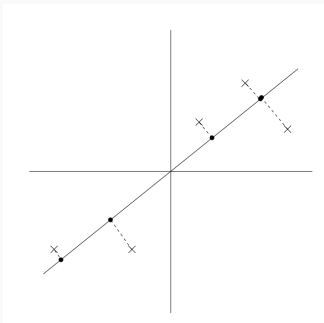


Рис. 6: Двумерный датасет



a)



b)

Рис. 7: Неверное и верное главное направление

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;
- обсудили метод PCA для предобработки данных.