

Методы “регуляризации” нейронных сетей

Яворская Мария

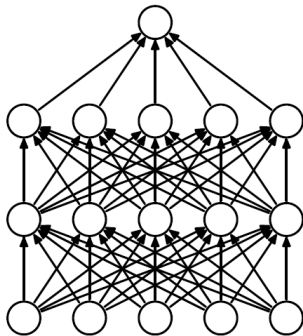
МГУ имени М. В. Ломоносова, факультет ВМК, кафедра ММП

22 марта 2018 г.

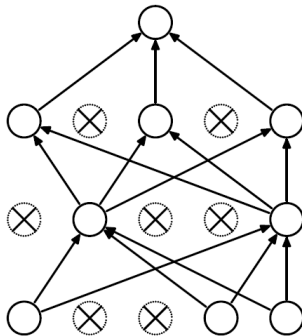
- Dropout, inverted dropout;
- Батч-нормализация;
- Методы инициализации.

- Проблемы:
 - переобучение — проблема, когда нейронная сеть теряет способность к обобщению;
 - сложность комбинирования ответов больших нейронных сетей во время тестирования;
- Решение: dropout;
- Идея: сети для обучения получаются с помощью исключения из сети (dropping out) нейронов. “Исключение” нейрона означает, что при любых входных данных или параметрах он возвращает 0.

Dropout



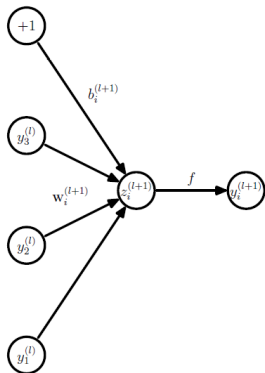
(a) Standard Neural Net



(b) After applying dropout.

Dropout

- $l \in \{1, \dots, L\}$ — индексы скрытых слоев сети;
- $\mathbf{z}^{(l)}$ — входной вектор слоя l ;
- $\mathbf{y}^{(l)}$ — выходной вектор слоя l ($\mathbf{y}^{(0)} = \mathbf{x}$ — исходные данные);
- $\mathbf{w}^{(l)}$ и $\mathbf{b}^{(l)}$ — вектор весов и сдвигов, соответствующих слою l ;
- $f()$ — функция активации.



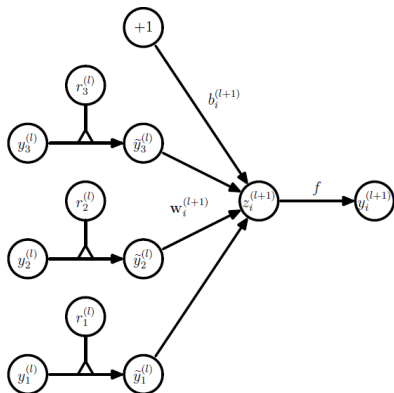
Прямой проход для стандартной нейронной сети:

$$z_i^{(l+1)} = w_i^{(l+1)} y^{(l)} + b_i^{(l+1)}$$

$$y_i^{(l+1)} = f(z_i^{(l+1)})$$

Dropout

- $l \in \{1, \dots, L\}$ — индексы скрытых слоев сети;
- $\mathbf{z}^{(l)}$ — входной вектор слоя l ;
- $\mathbf{y}^{(l)}$ — выходной вектор слоя l ($\mathbf{y}^{(0)} = \mathbf{x}$ — исходные данные);
- $\mathbf{w}^{(l)}$ и $\mathbf{b}^{(l)}$ — вектор весов и сдвигов, соответствующих слою l ;
- $f()$ — функция активации.



Прямой проход для нейронной сети с dropout:

$$r_j^{(l)} \sim \text{Bernoulli}(p)$$

$$\tilde{\mathbf{y}}^{(l)} = \mathbf{r}^{(l)} * \mathbf{y}^{(l)}$$

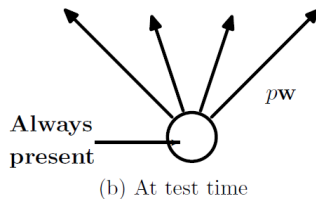
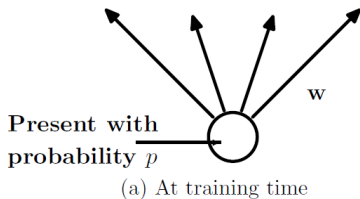
$$\mathbf{z}_i^{(l+1)} = \mathbf{w}_i^{(l+1)} \tilde{\mathbf{y}}^{(l)} + \mathbf{b}_i^{(l+1)}$$

$$\mathbf{y}_i^{(l+1)} = f(\mathbf{z}_i^{(l+1)})$$

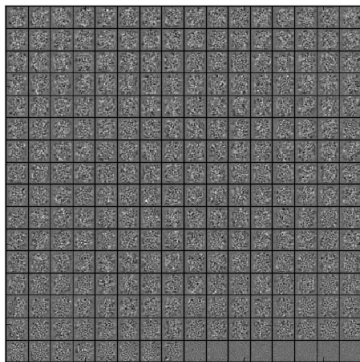
Dropout

Этап тестирования:

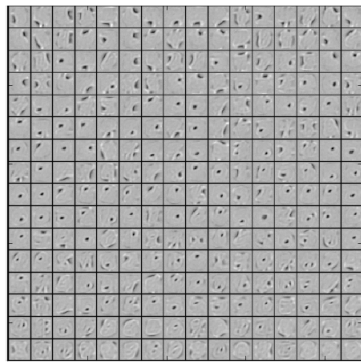
$$y_i^{(l+1)} = f(pw_i^{(l+1)}y^{(l)} + b_i^{(l+1)})$$



Dropout



(a) Without dropout



(b) Dropout with $p = 0.5$.

Признаки, выделенные на датасете MNIST при помощи автоэнкодера

Преимущества:

- Предотвращает переобучение;
- Возможность комбинирования ответов больших нейронных сетей во время тестирования.

Недостатки:

- Необходимость изменять нейронную сеть для проведения тестирования.

Inverted Dropout

Идея: умножать функцию активации на коэффициент не во время тестового этапа, а во время обучения.

Этап обучение:

$$y_i^{(l+1)} = f\left(\frac{1}{p} w_i^{(l+1)} (r^{(l)} * y^{(l)}) + b_i^{(l+1)}\right)$$

Этап тестирования:

$$y_i^{(l+1)} = f(w_i^{(l+1)} y^{(l)} + b_i^{(l+1)})$$

- Проблема: долгое время обучения нейронных сетей;
- Стандартные решения:
 - Увеличение learning rate;
 - Уменьшение число параметров сети;
 - Изменение learning rate в процессе обучения особым способом;
- Предлагаемое решение: батч-нормализация;
- Идея: добавить в архитектуру сети нормализацию по каждому из обучающих минибатчей.

Батч-нормализация

Дано: значения $\mathbf{x} \in B$, где $B = \{\mathbf{x}_{1\dots m}\}$

Требуется: $\{y_i = BN_{\gamma,\beta}(\mathbf{x}_i)\}$

Алгоритм:

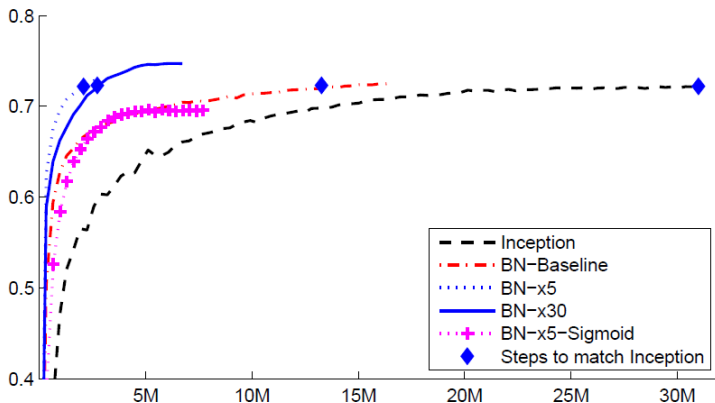
$$\textcircled{1} \quad \mu_B = \frac{1}{m} \sum_{i=1}^m \mathbf{x}_i$$

$$\textcircled{2} \quad \sigma_B^2 = \frac{1}{m} \sum_{i=1}^m (\mathbf{x}_i - \mu_B)^2$$

$$\textcircled{3} \quad \hat{\mathbf{x}}_i = \frac{\mathbf{x}_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

$$\textcircled{4} \quad y_i = \gamma \hat{\mathbf{x}}_i + \beta \equiv BN_{\gamma,\beta}(\mathbf{x}_i)$$

Батч-нормализация



Точность классификации «начального» алгоритма и его вариаций с батч-нормализацией vs количество итераций на этапе обучения.

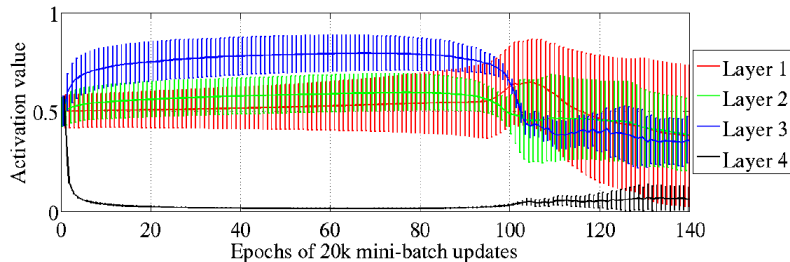
Преимущества:

- Уменьшение времени обучения нейронной сети;
- Возможность строить сети без применения dropout и с малым коэффициентом L_2 регуляризации.

Недостатки:

- Появляется еще $2L$ параметров, которые необходимо обучить.

- Проблема:
 - Перенасыщение нейронов;
 - Скорость сходимости.
- Решение: методы инициализации;
- Идея: осознанный выбор начальных значений весов для слоев, составляющих модель.



Средние значения и дисперсии выходных значений функций активаций различных слоев

Предположим, что веса и входные значения не коррелируют и имеют нулевое матожидание.

$$\begin{aligned} \text{Var}(\sum_{i=1}^{n_{in}} w_i x_i) &= \sum_{i=1}^{n_{in}} \text{Var}(w_i x_i) = \\ &= \sum_{i=1}^{n_{in}} \text{Var}(W) \text{Var}(X) = n_{in} \text{Var}(W) \text{Var}(X) \end{aligned}$$

Прямой проход:

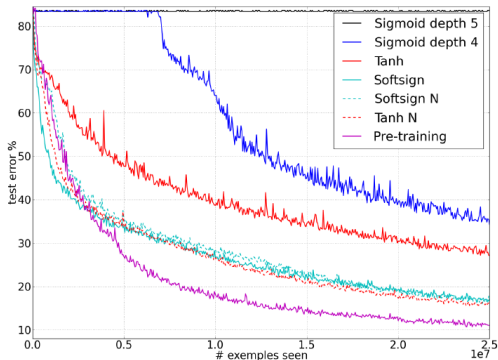
$$\text{Var}(W) = \frac{1}{n_{in}}$$

Обратный проход:

$$\text{Var}(W) = \frac{1}{n_{out}}$$

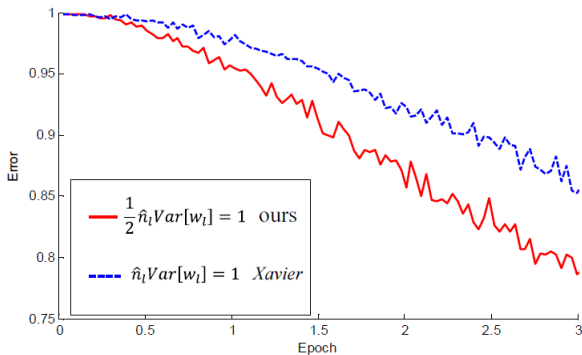
Метод инициализации Завьера (Xavier):

$$\text{Var}(W) = \frac{2}{n_{in} + n_{out}}$$



Метод инициализации Ге (He):

$$\text{Var}(W) = \frac{2}{n_{in}}$$



- Dropout:
 - Dropout: A Simple Way to Prevent Neural Networks from Overfitting — N.Srivastava, G.Hinton (2014)
 - Dropout — метод решения проблемы переобучения в нейронных сетях
- Батч-нормализация:
 - Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift — S.Ioffe, C.Szegedy (2015)
 - Batch Normalization для ускорения обучения нейронных сетей
- Методы инициализации:
 - Understanding the difficulty of training deep feedforward neural networks — X.Glorot, Y.Bengio (2010)
 - Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification — K.He, X.Zhang (2015)
 - An Explanation of Xavier Initialization