

Методы регуляризации нейронных сетей. Часть 2.

Weight Decay

Добавление к основной целевой функции

$$L_{reg} = \frac{\lambda}{2} \sum_i |w_i|^q,$$

$$q = 2 \text{ или } q = 1$$

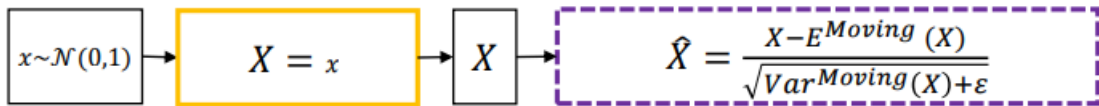
Dropout

Проблема (<https://arxiv.org/pdf/1801.05134.pdf>)



Train Mode $\text{Var}^{\text{Train}}(X) = \frac{1}{p} \rightarrow \text{Var}^{\text{Moving}}(X) = E\left(\frac{1}{p}\right)$

Test Mode $\text{Var}^{\text{Test}}(X) = 1 \not\rightarrow \text{Var}^{\text{Moving}}(X) = E\left(\frac{1}{p}\right)$



Dropout

Решение

1. Dropout после BatchNorm
2. Uout: добавление случайной величины с равномерным распределением к выходам

Dropout - Uout

$$X = x_i + x_i r_i$$

$$r_i \sim \mathcal{U}(-\beta, \beta)$$

$$0 \leq \beta \leq 1$$

Dropout - Uout

distributed random variable with zero mean and standard deviation equal to the activation of the unit, i.e., $x_i + x_i r$ and $r \sim \mathcal{N}(0, 1)$, we modify r as a uniform distribution that lies in $[-\beta, \beta]$, where $0 \leq \beta \leq 1$. Therefore, each hidden activation would be $X = x_i + x_i r_i$ and $r_i \sim \mathcal{U}(-\beta, \beta)$. We name this form of Dropout as “Uout” for simplicity. With the mutually independent distribution between x_i and r_i being hold, we apply the form $X = x_i + x_i r_i$, $r_i \sim \mathcal{U}(-\beta, \beta)$ in train stage and $X = x_i$ in test mode. Similarly, in the simplified case of $c = 0$, we can deduce the variance shift again as follows:

$$\begin{aligned}\frac{Var^{Test}(X)}{Var^{Train}(X)} &= \frac{Var(x_i)}{Var(x_i + x_i r_i)} = \frac{v}{E((x_i + x_i r_i)^2)} \\ &= \frac{v}{E(x_i^2) + 2E(x_i^2)E(r_i) + E(x_i^2)E(r_i^2)} = \frac{3}{3 + \beta^2}.\end{aligned}\tag{15}$$

Giving β as 0.1, the new variance shift rate would be $\frac{300}{301} \approx 0.9966777$ which is much closer to 1.0 than the previous

Label Smoothing

Используется и для борьбы с неправильной аннотацией.

$$E = - \sum_i \sum_k t_k^i \log(y_k(x_i))$$

One-hot encoding: $t_i = (0, 0, 1, 0)$, $t_i^k = 1$, $k = 2$ <- единица там, где правильный класс для данного наблюдения.

Меняем One-hot encoding вектор:

$$t_i = \left(\frac{\epsilon}{C-1}, \frac{\epsilon}{C-1}, 1 - \epsilon, \frac{\epsilon}{C-1} \right),$$

где C -- кол-во классов

Эвристики проверки нейронной сети

Начните с данных

Посмотрите, что

1. Они нормализованы
2. Категориальным переменным соответствует верное число класса, которыми они кодируются
3. Лейблы корректны для каждого наблюдения
4. Если есть предобработка данных, либо какая-то специфичная функция формирования батча, проконтролируйте, что лейблы по-прежнему правильные для каждого наблюдения
5. Выводите ошибки по каждому батчу, например, если это картинка, то каким картинкам присваиваются неверные классы.
6. Продебажьте полный pipeline, что на каждом шаге получается с данными: считывание, предобработка, подача в сеть.

Зафиксируйте **random seed** (везде!)

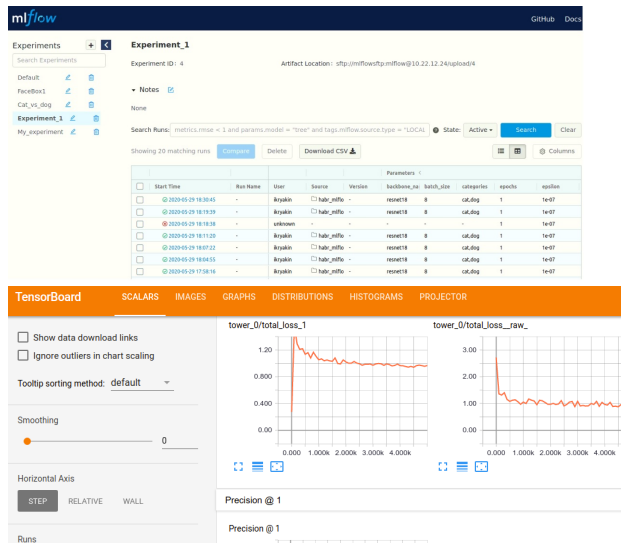
```
torch.manual seed(o)  
np.random.seed(o)
```

Начинайте с простого

Начинайте с простой предобработки данных.
Только после получения логичных результатов начинайте усложнение.

Логируйте все, что можно

Используйте TensorBoard, MLflow.



$\log(1 / \text{кол-во классов})$

Проверьте, что в самом начале loss выдает значение
 $\log(1 / \text{кол-во классов})$
для корректного класса.

Если не так, то скорее всего надо исправить инициализацию весов, в основном, править надо самый первый или последний слой.

Инициализация весов

Проверьте, что именно стоит по умолчанию.

Опираясь на распределение данных, возможно, стоит поправить среднее
или стандартное отклонение.

Также не забывайте про баесы, не всегда нужно инициализировать их
нулями.

Посмотрите, что обученная модель лучше **случайно**
проинициализированной модели.
Либо лучше модели, которая на вход принимает одни нули.

Заверфитьте на одном батче

Проверьте, что значение целевой функции сильно близко к нулю, а
точность под 100 %

Don't be a hero!

Выбирайте такую архитектуру, которая уже дала хороший результат.
Берите то, что пишут в статьях, но не сложное.
Усложняйте модель постепенно.

Используйте советы, которые дают другие.
Берите к сведению любые эвристики.
Например, что начинайте с Adam и ReLU.

Нет понятия маленького или большого Learning Rate.

Есть маленький и большой Learning Rate для Вашей (!) задачи.

Обращайте на его начальное значение.

Вначале делайте его константным, только после получения логичных результатов вводите политику его изменения.

Не переусердствуйте с **Dropout**, **BatchNorm** и **Weight Decay**.

Просто дайте нейронной сети **время**.
Не всегда сеть может обучиться быстро.