Нейронные сети: как заставить их обучаться

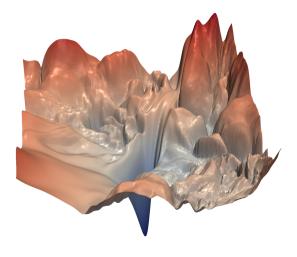
Константин Архипенко

7 октября 2020 г.



Почему сложно обучать нейросети



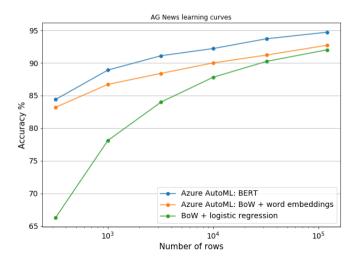


Почему сложно обучать нейросети



- Много "плохих" локальных минимумов и седловых точек
- Переобучение
- · Transfer learning полностью не спасает





Аугментация данных



Fig. 1. Exemplar applications of image transformations available in Albumentations

• Для картинок: повороты, размытие, масштабирование, сгенерированные GAN, ...

²Buslaev et al. "Albumentations: Fast and Flexible Image Augmentations" (2020)

Аугментация для NLP



Augmenter	Target	Augmenter	Action	Description
Textual	Character	KeyboardAug	substitute	Simulate keyboard distance error
Textual		OcrAug	substitute	Simulate OCR engine error
Textual		RandomAug	insert, substitute, swap, delete	Apply augmentation randomly
Textual	Word	AntonymAug	substitute	Substitute opposite meaning word according to WordNet antonym
Textual		Contextual Word Embs Aug	insert, substitute	Feeding surroundings word to BERT, DistilBERT, ROBERTa or XLNet language model to find out the most suitlabe word for augmentation



	Sentence	
Original	The quick brown fox jumps over the lazy dog	
Synonym (PPDB)	The quick brown fox climbs over the lazy dog	
Word Embeddings (word2vec)	The easy brown fox jumps over the lazy dog	
Contextual Word Embeddings (BERT)	Little quick brown fox jumps over the lazy dog	
PPDB + word2vec + BERT	Little easy brown fox climbs over the lazy dog	

³https://github.com/makcedward/nlpaug



- Пример логистическая регрессия (бинарная классификация)
- · Maximum likelihood estimation:

$$\frac{1}{N} \sum_{n} \log P(y = y_n | \mathbf{x}_n, \mathbf{w}, w_0) \to \max_{\mathbf{w}, w_0}$$

$$\hat{y} = P(y = 1 | \mathbf{x}, \mathbf{w}, w_0) = \sigma(w_0 + \langle \mathbf{w}, \mathbf{x} \rangle)$$

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

• Переход к минимизации кросс-энтропии:

$$\frac{1}{N} \sum_{n} \operatorname{xent}(y_n, \hat{y}_n(\mathbf{x}_n, \mathbf{w}, w_0)) \to \min_{\mathbf{w}, w_0}$$

$$\operatorname{xent}(y, \hat{y}) = -(y \log \hat{y} + (1 - y) \log(1 - \hat{y}))$$



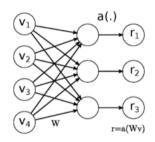
- Уравнение разделяющей гиперплоскости: $w_0 + \langle \mathbf{w}, \mathbf{x} \rangle = 0$
- Умножение уравнения на 2 не меняет саму гиперплоскость
- Но увеличение w делает модель более чувствительной к выбросам
- Логистическая регрессия с ℓ_2 -регуляризацией:

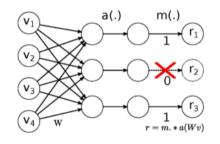
$$\frac{1}{N}\sum_{n}\operatorname{xent}(y_{n},\hat{y}_{n}(\mathbf{x}_{n},\mathbf{w},w_{0}))+\frac{\lambda}{2}\|\mathbf{w}\|_{2}\to\min_{\mathbf{w},w_{0}}$$

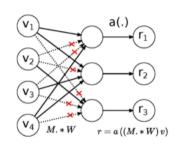
• В нейронных сетях: регуляризация dense слоев и сверточных фильтров

Dropout и DropConnect









No-Drop Network

DropOut Network

DropConnect Network



- \cdot p drop rate, $a(\cdot)$ функция активации, **W** матрица параметров
- · Dropout:

$$m_i \sim \text{Bernoulli}(p)$$

$$\mathbf{r} = \frac{1}{1 - \rho} \mathbf{m} \odot a(\mathbf{W} \mathbf{v})$$

DropConnect:

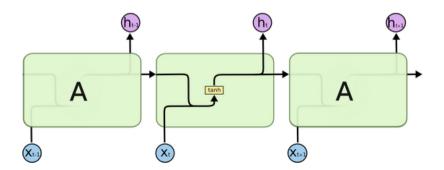
$$m_{ij} \sim \text{Bernoulli}(p)$$

$$r = a((M \odot W)v)$$



- · До появления Transformer SOTA языковые модели
- · Vanilla RNN:

$$h_t = f(Wx_t + Uh_{t-1} + b)$$
$$h_0 = 0$$





$$\mathbf{h}_t = f(\mathbf{W}\mathbf{x}_t + \mathbf{U}\mathbf{h}_{t-1} + \mathbf{b})$$

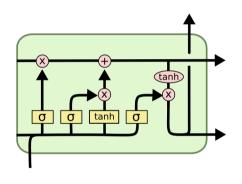
• Проблемы с градиентом:

$$\frac{\partial \mathcal{L}}{\partial W} = \frac{\partial \mathcal{L}}{\partial h_T} \frac{\partial h_T}{\partial W} \qquad \frac{\partial h_T}{\partial W} = (\prod_{t=1}^{l-1} \frac{\partial h_{t+1}}{\partial h_t}) \cdot \frac{\partial h_1}{\partial W}$$

• Величина (норма) произведения зависит от нормы матрицы U. Если последняя слишком отклоняется от 1, получим vanishing/exploding gradients

Long short-term memory





$$\begin{split} \mathbf{i}_t &= \sigma(\mathbf{W}^i \mathbf{x}_t + \mathbf{U}^i \mathbf{h}_{t-1} + \mathbf{b}^i) \\ \mathbf{f}_t &= \sigma(\mathbf{W}^f \mathbf{x}_t + \mathbf{U}^f \mathbf{h}_{t-1} + \mathbf{b}^f) \\ \mathbf{o}_t &= \sigma(\mathbf{W}^o \mathbf{x}_t + \mathbf{U}^o \mathbf{h}_{t-1} + \mathbf{b}^o) \\ \widetilde{\mathbf{c}}_t &= \tanh(\mathbf{W}^c \mathbf{x}_t + \mathbf{U}^c \mathbf{h}_{t-1} + \mathbf{b}^c) \\ \mathbf{c}_t &= \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \widetilde{\mathbf{c}}_t \\ \mathbf{h}_t &= \tanh(\mathbf{c}_t) \odot \mathbf{o}_t \end{split}$$



- Merity et al. "Regularizing and Optimizing LSTM Language Models" (2017)
- Применяется DropConnect к рекуррентным матрицам ($\mathbf{U}^i,\mathbf{U}^f,\mathbf{U}^c,\mathbf{U}^o$)
- Sample weight отношение длины обрезанной последовательности к длине исходной
- Weight tying
- Activity regularization



- Weight tying: используются одни и те же веса для embedding и softmax слоев
- Activity regularization: ℓ_2 -регуляризация скрытых состояний \mathbf{h}_t
- Temporal activity regularization: ℓ_2 -регуляризация разницы между \mathbf{h} для соседних timestep:

$$TAR = \beta \|\mathbf{h}_t - \mathbf{h}_{t-1}\|_2$$

- Вывод: SOTA модели часто содержат специфичные трюки, касающихся регуляризации и оптимизации
- Сейчас рекуррентные нейросети для моделирования языка вытеснены Трансформерами

Label smoothing⁴

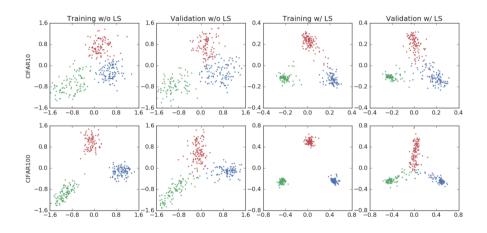


- Классификация на К классов
- Ground truth one-hot векторы \mathbf{y}_n с единицей в позиции номера класса
- \cdot Заменяем единицу в этих векторах на 1 lpha, а все нули на $rac{lpha}{K-1}$

⁴Müller et al. "When Does Label Smoothing Help?" (2019)

Label smoothing



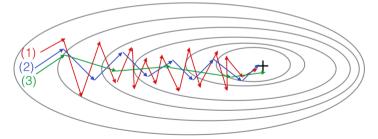




· Stochastic gradient descent c momentum:

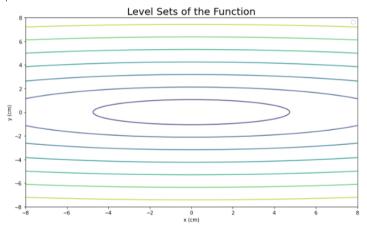
$$\mathbf{v}_{t} \leftarrow \gamma \mathbf{v}_{t-1} - \eta \nabla_{\boldsymbol{\theta}} L(B, \boldsymbol{\theta})|_{\boldsymbol{\theta} = \boldsymbol{\theta}_{t-1}, B = B_{t}}$$
$$\boldsymbol{\theta}_{t} \leftarrow \boldsymbol{\theta}_{t-1} + \mathbf{v}_{t}$$

• Типичное значение $\gamma = 0.9$





• Мотивация:





• Скользящее среднее квадрата градиента:

$$Eg_t^2 = \beta Eg_{t-1}^2 + (1 - \beta)g_t^2$$
$$\boldsymbol{\theta}_t \leftarrow \boldsymbol{\theta}_{t-1} - \frac{\eta}{\sqrt{Eg_t^2 + \varepsilon}} \cdot g_t$$

- Вторая производная связь с кривизной по направлению
- Adagrad: сумма вместо скользящего среднего, эффективная learning rate быстро падает (недостаток)



$$g_t \leftarrow \nabla_{\boldsymbol{\theta}} L(B, \boldsymbol{\theta})|_{\boldsymbol{\theta} = \boldsymbol{\theta}_{t-1}, B = B_t}$$

$$\widetilde{\boldsymbol{\phi}}_{t} \leftarrow \beta_{1} \widetilde{\boldsymbol{\phi}}_{t-1} + (1 - \beta_{1}) \cdot \mathbf{g}_{t}$$

$$\widetilde{\boldsymbol{\psi}}_{t} \leftarrow \beta_{2} \widetilde{\boldsymbol{\psi}}_{t-1} + (1 - \beta_{2}) \cdot \mathbf{g}_{t}^{2}$$

$$\boldsymbol{\phi}_t \leftarrow \widetilde{\boldsymbol{\phi}}_t / (1 - \beta_1^t)$$

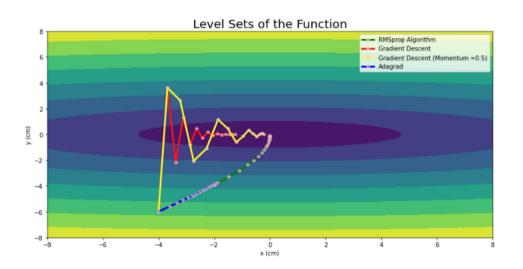
$$\boldsymbol{\psi}_t \leftarrow \widetilde{\boldsymbol{\psi}}_t / (1 - \beta_2^t)$$

$$\boldsymbol{\theta}_t \leftarrow \boldsymbol{\theta}_{t-1} - \alpha \boldsymbol{\phi}_t / (\sqrt{\boldsymbol{\psi}_t} + \varepsilon)$$

- Скользящие средние градиента и квадрата градиента
- Проблема холодного старта: деление на $(1-\beta)^t$

⁵Kingma & Ba "Adam: A Method for Stochastic Optimization" (2014) ICLR







- SGD c momentum обычно требуется больше времени для сходимости, но он менее подвержен over-shooting (как Adam)
- Сказанное выше типичное поведение, не всегда оно так
- Есть статья, где предлагается динамическое переключение с Adam на SGD^6

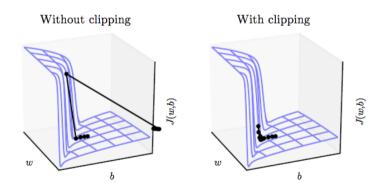
⁶Keskar & Socher "Improving Generalization Performance by Switching from Adam to SGD" (2017)

Другие трюки для оптимизации



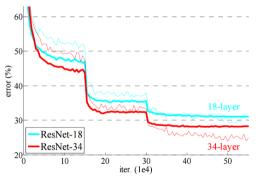
• Gradient clipping: если норма градиента больше с, то

$$\mathbf{g} \leftarrow \mathbf{c} \cdot \frac{\mathbf{g}}{\|\mathbf{g}\|}$$





· Learning rate schedule:



- \cdot Похожий эффект имеет изменение batch size во время обучения 7
- · Early stopping

⁷Smith et al. "Don't Decay the Learning Rate, Increase the Batch Size" ICLR 2018

Батч-нормализация

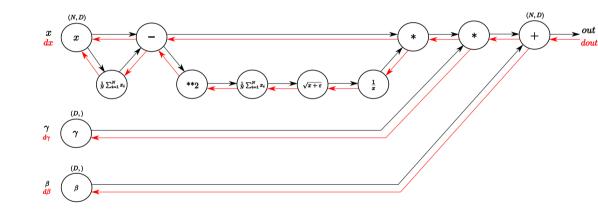


```
Input: Values of x over a mini-batch: \mathcal{B} = \{x_{1...m}\};
               Parameters to be learned: \gamma, \beta
Output: \{y_i = BN_{\gamma,\beta}(x_i)\}
  \mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^{m} x_i
                                                                        // mini-batch mean
  \sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2
                                                           // mini-batch variance
   \widehat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}
                                                                                     // normalize
     u_i \leftarrow \gamma \widehat{x}_i + \beta \equiv BN_{\gamma,\beta}(x_i)
                                                                            // scale and shift
```

• Batch norm после функции активации часто дает результаты лучше_{27/30}

Batch norm: граф вычислений





Почему batch norm работает?



- Раньше считалось, что помогает справиться с internal covariate shift⁸ (постоянное изменение распределения входа следующего слоя при изменении весов предыдущего)
- Более новый результат: делает функцию потерь более гладкой⁹

Theorem 4.1 (The effect of BatchNorm on the Lipschitzness of the loss). *For a BatchNorm network with loss* $\widehat{\mathcal{L}}$ *and an identical non-BN network with (identical) loss* \mathcal{L} ,

$$\left|\left|\nabla_{\boldsymbol{y_j}}\widehat{\mathcal{L}}\right|\right|^2 \leq \frac{\gamma^2}{\sigma_j^2} \left(\left|\left|\nabla_{\boldsymbol{y_j}}\mathcal{L}\right|\right|^2 - \frac{1}{m} \left\langle \mathbf{1}, \nabla_{\boldsymbol{y_j}}\mathcal{L}\right\rangle^2 - \frac{1}{m} \left\langle \nabla_{\boldsymbol{y_j}}\mathcal{L}, \hat{\boldsymbol{y}}_j\right\rangle^2\right).$$

⁸Ioffe & Szegedy "Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift" (2015)

⁹Santurkar et al. "How Does Batch Normalization Help Optimization? " (2018)

Зависимость BN-сети от инициализации



Lemma 4.5 (BatchNorm leads to a favourable initialization). Let W^* and \widehat{W}^* be the set of local optima for the weights in the normal and BN networks, respectively. For any initialization W_0

$$\left| \left| W_0 - \widehat{W}^* \right| \right|^2 \le \left| \left| W_0 - W^* \right| \right|^2 - \frac{1}{\left| \left| W^* \right| \right|^2} \left(\left| \left| W^* \right| \right|^2 - \langle W^*, W_0 \rangle \right)^2,$$

if $\langle W_0, W^* \rangle > 0$, where \widehat{W}^* and W^* are closest optima for BN and standard network, respectively.

- Доказательство: $BN((a\mathbf{W})\mathbf{x}) = BN(\mathbf{W}\mathbf{x}) \ \, \forall a>0$, поэтому для любого минимума в non-BN сети \mathbf{W}^* точка $\hat{\mathbf{W}}=k\mathbf{W}^*$ будет минимумом в BN-сети
- Возьмем $k = \frac{\langle \mathbf{W}^*, \mathbf{W}_0 \rangle}{\|\mathbf{W}^*\|^2}$. Раскроем квадраты норм в $\|\mathbf{W}_0 \hat{\mathbf{W}}\|^2 \|\mathbf{W}_0 \mathbf{W}^*\|^2$, получим отрицательное число $-\|\mathbf{W}^*\|^2 \cdot (1-k)^2$
- Последнее верно для произвольного $\hat{\mathbf{W}}$, значит, и для $\hat{\mathbf{W}}^*$