Основы практического использования нейронных сетей.

Лекция 4. Методы повышения эффективности алгоритмов обучения для глубоких НС.

Дмитрий Буряк. к.ф.-м.н. dyb04@yandex.ru

Batch normalization

 \Box Нормализация + декорреляция входных данных \rightarrow повышение эффективности обучения. □ Нарушение полученных свойств входных векторов для промежуточных данных во внутренних слоях (internal covariance shift). □ Идея: проводить предобработку входных данных для каждого внутреннего слоя. $lue{}$ Оптимизация вычислительных затрат o нормализация внутренних данных (без декорреляции).

Sergey Ioffe, Christian Szegedy. Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift, 2015

Batch normalization. Алгоритм

```
Input: Network N with trainable parameters \Theta;
                             subset of activations \{x^{(k)}\}_{k=1}^{K}
Output: Batch-normalized network for inference, Normalized network for inference network network for inference network n
   1: N_{\text{BN}}^{\text{tr}} \leftarrow N // Training BN network
   2: for k = 1 ... K do
   3: Add transformation y^{(k)} = BN_{\gamma(k),\beta(k)}(x^{(k)}) to
                         N_{\rm RN}^{\rm tr} (Alg. 1)
   4: Modify each layer in N_{\text{RN}}^{\text{tr}} with input x^{(k)} to take
                         y^{(k)} instead
   5: end for
   6: Train N_{\rm BN}^{\rm tr} to optimize the parameters \Theta \cup
               \{\gamma^{(k)}, \beta^{(k)}\}_{k=1}^{K}
   7: N_{\rm BN}^{\rm inf} \leftarrow N_{\rm BN}^{\rm tr} // Inference BN network with frozen
                                                                             // parameters
   8: for k = 1 ... K do
                      // For clarity, x \equiv x^{(k)}, \gamma \equiv \gamma^{(k)}, \mu_B \equiv \mu_B^{(k)}, etc.
                     Process multiple training mini-batches \mathcal{B}, each of
                         size m, and average over them:
                                                                                      E[x] \leftarrow E_{\mathcal{B}}[\mu_{\mathcal{B}}]
                                                                                Var[x] \leftarrow \frac{m}{m-1} E_B[\sigma_B^2]
                    In N_{\rm BN}^{\rm inf}, replace the transform y = BN_{\gamma,\beta}(x) with
                        y = \frac{\gamma}{\sqrt{\text{Var}[x] + \epsilon}} \cdot x + (\beta - \frac{\gamma E[x]}{\sqrt{\text{Var}[x] + \epsilon}})
12: end for
```

```
Input: Values of x over a mini-batch: \mathcal{B} = \{x_{1...m}\};

Parameters to be learned: \gamma, \beta

Output: \{y_i = \mathrm{BN}_{\gamma,\beta}(x_i)\}

\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \qquad // \text{mini-batch mean}
\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \qquad // \text{mini-batch variance}
\widehat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \qquad // \text{normalize}
y_i \leftarrow \gamma \widehat{x}_i + \beta \equiv \mathrm{BN}_{\gamma,\beta}(x_i) \qquad // \text{scale and shift}
```

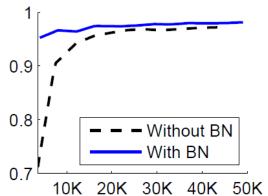
Batch normalization. Примеры.

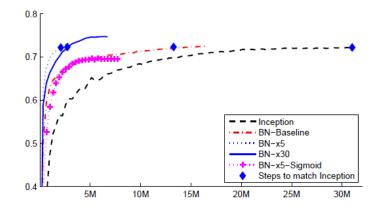
- Классификация MNIST
- □ Сеть 784x100x100x100x10,
- 50000 обучающих примеров

- Классификация ImageNet
- **□** Сеть 13.6*10⁶ параметров,

1000 классов

Model	Steps to 72.2%	Max accuracy
Inception	$31.0 \cdot 10^{6}$	72.2%
BN-Baseline	$13.3 \cdot 10^{6}$	72.7%
BN-x5	$2.1 \cdot 10^{6}$	73.0%
BN-x30	$2.7 \cdot 10^{6}$	74.8%
BN-x5-Sigmoid		69.8%





Регуляризация L2

- □ Регуляризация метод предотвращения переобучения НС.
- □ Введение штрафа для больших весов.

$$C = C_0 + rac{\lambda}{2n} \sum_w w^2$$

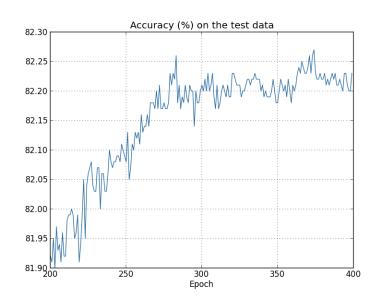
🗖 λ - коэффициент регуляризации.

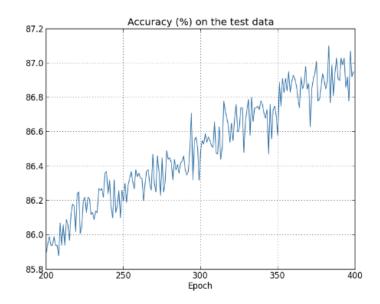
$$egin{aligned} rac{\partial C}{\partial w} &= rac{\partial C_0}{\partial w} + rac{\lambda}{n} w \ rac{\partial C}{\partial b} &= rac{\partial C_0}{\partial b}. \ rac{\partial C}{\partial b} &= rac{\partial C_0}{\partial b}. \end{aligned} \qquad egin{aligned} b &
ightarrow b - \eta rac{\partial C_0}{\partial b}. \ w &
ightarrow w - \eta rac{\partial C_0}{\partial w} - rac{\eta \lambda}{n} w \ &= \left(1 - rac{\eta \lambda}{n}
ight) w - \eta rac{\partial C_0}{\partial w}. \end{aligned}$$

□ Масштабирование веса перед коррекцией по градиентному спуску.

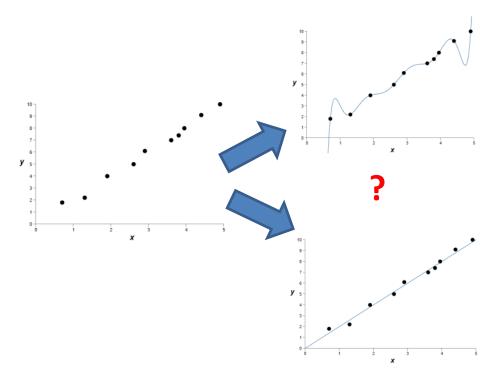
Пример применения регуляризации L2

- Классификация MNIST
- □ Сеть 784х30х10, 1000 обучающих примеров





Регуляризация → снижение переобучения



- ☐ Нет однозначного решения без дополнительной инфоормации.
- □ Большие значения параметров → увеличение чувствительности к шуму.

$$y = a_0 x^9 + a_1 x^8 + \dots$$

$$y = a_0 x + a_1$$

Регуляризация L1

□ Введение штрафа для больших весов.

$$C = C_0 + \frac{\lambda}{n} \sum_{w} |w|$$

🗖 λ - коэффициент регуляризации.

$$rac{\partial C}{\partial w} = rac{\partial C_0}{\partial w} + rac{\lambda}{n} \operatorname{sgn}(w) \qquad w o w' = w - rac{\eta \lambda}{n} \operatorname{sgn}(w) - \eta rac{\partial C_0}{\partial w}$$

- □ Уменьшение веса на фиксированную величину
- □ Для регуляризации L2 значение уменьшения веса зависит от его величины.

$$w o w'=w\left(1-rac{\eta\lambda}{n}
ight)-\etarac{\partial C_0}{\partial w}$$

Регуляризация Max-norm

□ Ограничения нормы вектора весов для каждого нейрона.

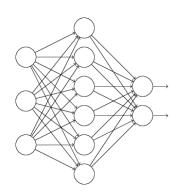
$$||\mathbf{w}||_2 \le c$$

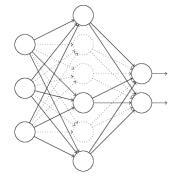
- **а** с гиперпараметр.
- □ Реализация через нормировку вектора весов при
- невыполнении неравенства.
- □ Эффективна при совместном использовании с Dropout

Dropout

- □ Инструмент регуляризации.
- □ Модификация архитектуры сети в процессе обучения.
- Упрощенная схема Dropout
 - 1. Временно удалить из НС половину случайно выбранных внутренних нейронов с соответствующими связями.
 - 2. Провести итерацию обучение на пакете: обновление связей оставшихся нейронов.
 - 3. Восстановить удаленные нейроны и их связи.
 - 4. Повторить π . 1 3.
- □ Перед применением сети уменьшить внутренние веса в 2 раза.

N. Srivastava et al. Dropout: A simple way to prevent neural networks from overfitting, 2014.





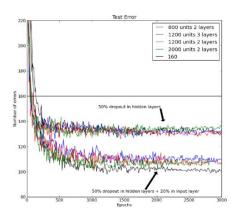
Dropout. Пример использования.

☐ MNIST.

Входной вектор 784 элемента.

10 классов.

10000 тестовых изображений.

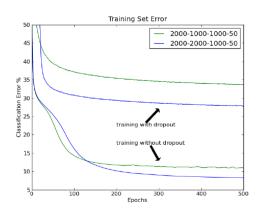


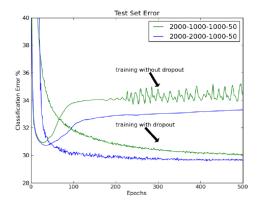
☐ Reuters. Классификация документов.

Входной вектор 2000 элементов.

50 классов.

~200000 тестовых документов.

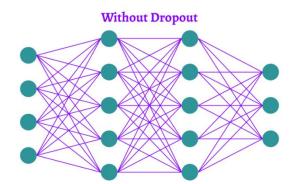


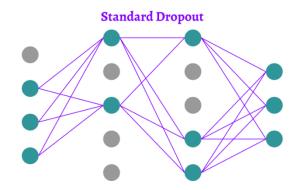


G.E. Hinton et al, Improving neural networks by preventing co-adaptation of feature detectors, 2012

Варианты Dropout. DropConnect.

□ Удаление связей.



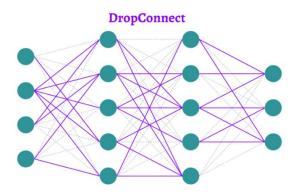


Training Phase:

$$y = f(Wx) \circ m$$
, $m_i \sim Bernoulli(p)$

Testing Phase:

$$y = (1 - p)f(Wx)$$



Training Phase:

$$\mathbf{y} = f((\mathbf{W} \circ \mathbf{M})\mathbf{x}), \quad M_{i,j} \sim Bernoulli(p)$$

Testing Phase:

$$y = (Wx) \circ \hat{m}(Z)$$

where
$$\hat{m}_i(Z) = \frac{1}{Z} \sum_{z=0}^{Z} f(\hat{x}_{i,z}), \quad \hat{x}_{i,z} \sim \mathcal{N}(\mu_i, \sigma_i^2)$$

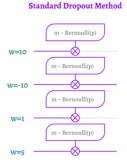
and
$$\boldsymbol{\mu} = p\mathbf{W}\mathbf{x}$$
, $\boldsymbol{\sigma}^2 = p(1-p)(\mathbf{W} \circ \mathbf{W})(\mathbf{x} \circ \mathbf{x})$, $Z \in \mathbb{N}^+$

Варианты Dropout. Standout.

- Вероятность удаления нейрона зависит от величины весов.
- □ Пример.

$$g(x) = |\sigma(x)|$$

$$W_s = \alpha W + \beta$$



Training Phase:

$$y = f(Wx) \circ m$$
, $m_i \sim Bernoulli(g(W_sx))$

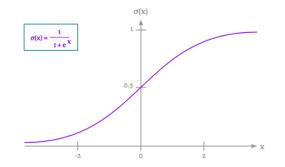
Testing Phase:

$$y = (1 - g(W_sx)) \circ f(Wx)$$

where W_s is the belief network's weights and g is the belief network's activation function







Варианты Dropout. Gaussian Dropout.

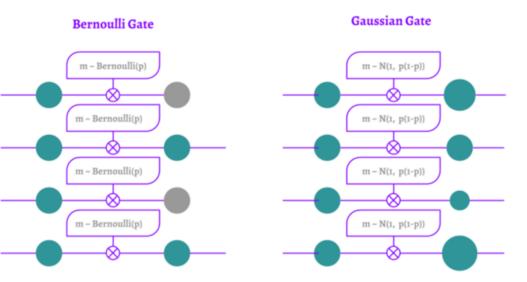
- □ Нейроны не удаляются, а «завешиваются» с помощью нормального распределения.
- □ Выше скорость сходимости.

Training Phase:

$$y = f(\mathbf{W}\mathbf{x}) \circ \mathbf{m}, \quad m_i \sim \mathcal{N}(1, p(1-p))$$

Testing Phase:

$$y = f(Wx)$$



Варианты Dropout. Pooling Dropout.

☐ Стандартный dropout не Without Max-Pooling Dropout

With Max-Pooling Dropout

эффективен для

изображений

☐ Pooling Dropout

применяют для сверточных

сетей

Training Phase:

$$Y = max\{Pool_{size}(Y) \circ M_{size}\}$$
 $M_{ij} \sim Bernoulli(p)$

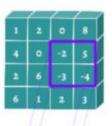
Testing Phase:

$$\mathbf{Y} = (1 - p) \max\{Pool_{size}(\mathbf{Y})\}$$







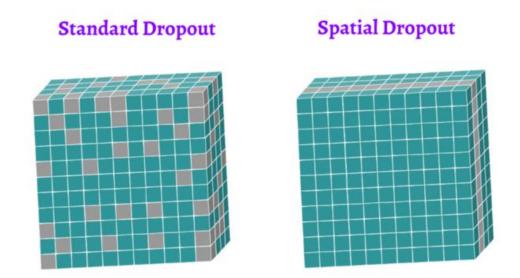






Варианты Dropout. Spatial Dropout.

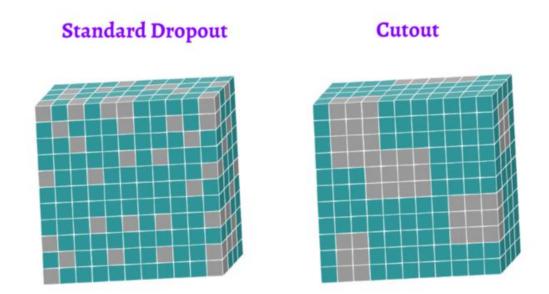
- □ Удаление карт признаков
- ☐ Spatial Dropout применяют для сверточных сетей



https://towardsdatascience.com/12-main-dropout-methods-mathematical-and-visual-explanation-58cdc2112293

Варианты Dropout. Cutout.

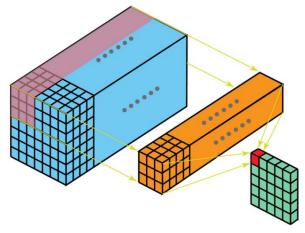
- □ Удаление фрагментов карт
- ☐ Cutout применяют для сверточных сетей



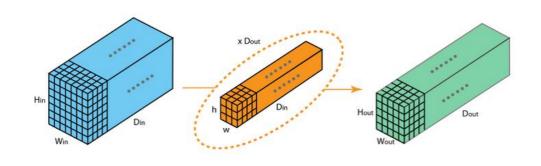
https://towardsdatascience.com/12-main-dropout-methods-mathematical-and-visual-explanation-58cdc2112293

Типы сверток (1)

□ 2D свертка



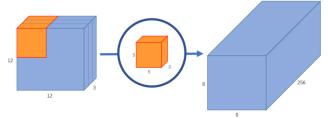
2D свертка: получение одного выходного элемента



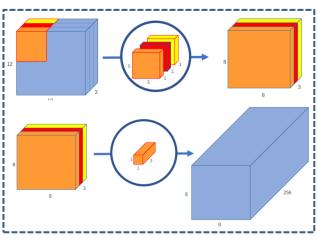
Выполнение 2D свертки, число входных карт Din, выходных - Dout

Типы сверток (2)

- □ Входной тензор: WxHxC_i;
- \square Ядро: $K_1 \times K_2$;
- \square Выходной тензор: WxHxC $_{0}$;
- Обычная свертка
 - Число операций: $K_1xK_2xC_1xWxHxC_0$;
- Свертка Depthwise Separable
 - Depthwise, число операций: K₁xK₂xC_ixWxH;
 - Pointwise, число операций: C_IxWxHxC_O;
 - Общее число операций: K₁xK₂xC₁xWxH+ C₁xWxHxC₀;
 - +: Вычислительная эффективность
 - -: Потеря точности



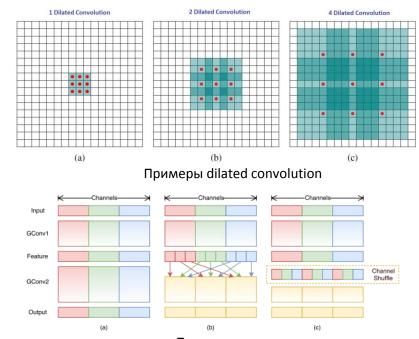
Обычная свертка, 1228800 операций



Depthwise Separable, 52952 операций

Типы сверток (3)

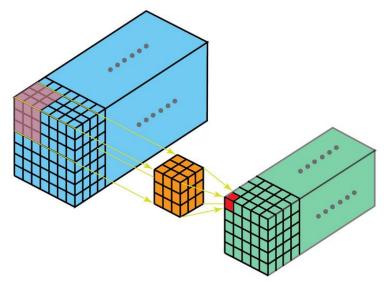
- ☐ Dilated convolution
 - Прореживание рецептивного поля
 - Увеличение рецептивного поля с сохранением числа параметров и операций.
- ☐ Shuffle convolution
 - Поканальные свертки с дальнейшим перемешиванием результирующих карт.



a. Поканальная свертка b, c. Shuffle convolution

Типы сверток (4)

- □ 3D свертка
- □ Применяют для обработки 3D данных



Выполнение 3D свертки

Вопросы

- □ Почему Batch Normalization повышает эффективность обучения?
- ☐ Какие типы Dropout применяют для сверточных сетей?
- □ Почему Depthwise Separable свертка является вычислительно более

эффективной во сравнению с обычной операцией свертки в НС?