Основы практического использования нейронных сетей.

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

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



Пример НС для распознавания изображений

- □ Распознавание изображений кошек и собак (база cats vs dogs)
 - RGB, различное разрешение, 25000 изображений.

```
→ inputs = keras.Input(shape=(180, 180, 3))
 The model
             x = layers.Rescaling(1./255)(inputs)
   expects
             x = layers.Conv2D(filters=32, kernel size=3, activation="relu")(x)
RGB images
             x = layers.MaxPooling2D(pool size=2)(x)
    of size
             x = layers.Conv2D(filters=64, kernel_size=3, activation="relu")(x)
180 \times 180.
             x = layers.MaxPooling2D(pool size=2)(x)
             x = layers.Conv2D(filters=128, kernel size=3, activation="relu")(x)
             x = layers.MaxPooling2D(pool size=2)(x)
             x = layers.Conv2D(filters=256, kernel size=3, activation="relu")(x)
             x = layers.MaxPooling2D(pool size=2)(x)
             x = layers.Conv2D(filters=256, kernel size=3, activation="relu")(x)
             x = layers.Flatten()(x)
             outputs = layers.Dense(1, activation="sigmoid")(x)
             model = keras.Model(inputs=inputs, outputs=outputs)
```

Rescale inputs to the [0, 1] range by dividing them by 255.

Подготовка обучающих и тестовых данных

☐ Использование объекта Dataset.

```
from tensorflow.keras.utils import image_dataset_from_directory

train_dataset = image_dataset_from_directory(
    new_base_dir / "train",
    image_size=(180, 180),
    batch_size=32)

validation_dataset = image_dataset_from_directory(
    new_base_dir / "validation",
    image_size=(180, 180),
    batch_size=32)

test_dataset = image_dataset_from_directory(
    new_base_dir / "test",
    image_size=(180, 180),
    batch_size=32)
```

```
history = model.fit(
    train_dataset,
    epochs=30,
    validation_data=validation_dataset,
    callbacks=callbacks)
```

Аугментация данных

- □ Формирование синтетических данных для повышения разнообразия данных при обучении.
 - Необходимо учитывать, какие реальные искажения данных возможны/допустимы















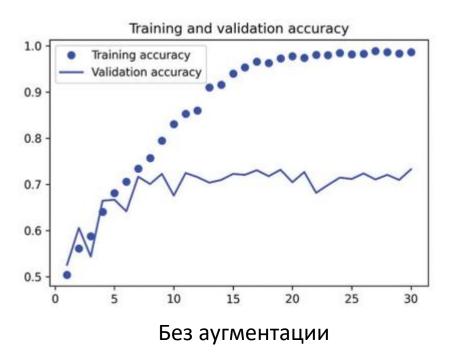


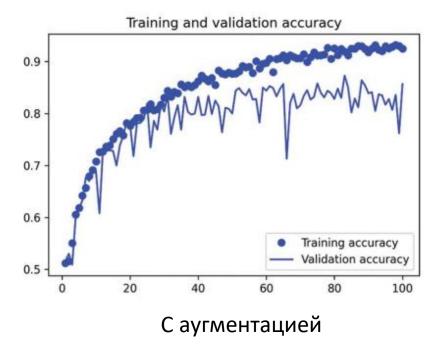


НС с аугментацией данных

```
inputs = keras.Input(shape=(180, 180, 3))
x = data_augmentation(inputs)
x = layers.Rescaling(1./255)(x)
x = layers.Conv2D(filters=32, kernel size=3, activation="relu")(x)
x = layers.MaxPooling2D(pool_size=2)(x)
x = layers.Conv2D(filters=64, kernel_size=3, activation="relu")(x)
x = layers.MaxPooling2D(pool size=2)(x)
x = layers.Conv2D(filters=128, kernel size=3, activation="relu")(x)
x = layers.MaxPooling2D(pool size=2)(x)
x = layers.Conv2D(filters=256, kernel size=3, activation="relu")(x)
x = layers.MaxPooling2D(pool size=2)(x)
x = layers.Conv2D(filters=256, kernel size=3, activation="relu")(x)
x = layers.Flatten()(x)
```

Влияние аугментации на точность на тестовых данных





Регуляризация L2

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

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

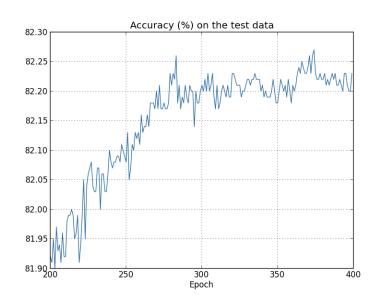
🗖 λ - коэффициент регуляризации.

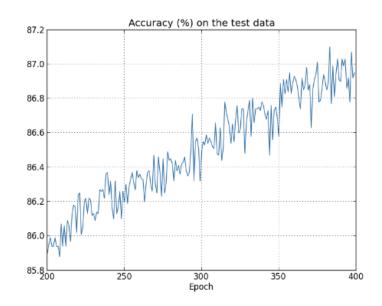
$$egin{aligned} rac{\partial C}{\partial w} &= rac{\partial C_0}{\partial w} + rac{\lambda}{n} w \ rac{\partial C}{\partial b} &= rac{\partial C_0}{\partial b}. \ rac{\partial C}{\partial b} &= rac{\partial C_0}{\partial b}. \end{aligned} \qquad egin{aligned} b &
ightarrow b - \eta rac{\partial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

□ Масштабирование веса перед коррекцией по градиентному спуску.

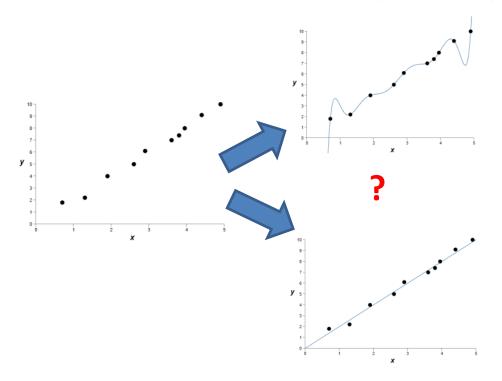
Пример применения регуляризации L2

- 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.

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

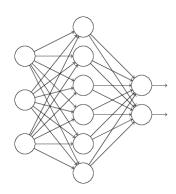
Регуляризация в Keras

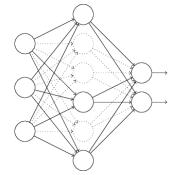
- □ keras.layers.Dense(units, ..., kernel_regularizer=None, bias_regularizer=None,)
 □ keras.layers.Conv2D(filters, kernel_size, ..., kernel_regularizer=None, bias_regularizer=None,)
 □ Доступные регуляризаторы:
 - keras.regularizers.L1(l1=0.01)
 - keras.regularizers.L2(l2=0.01)
 - keras.regularizers.L1L2(l1=0.01, l2=0.01)
 - keras.regularizers.OrthogonalRegularizer(factor=0.01, mode="rows")
 - реализация собственных регуляризаторов

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.





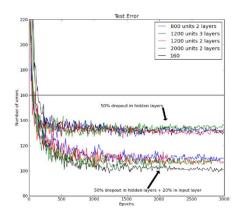
Dropout. Пример использования.

☐ MNIST.

Входной вектор 784 элемента.

10 классов.

10000 тестовых изображений.

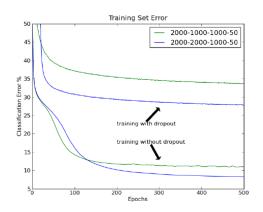


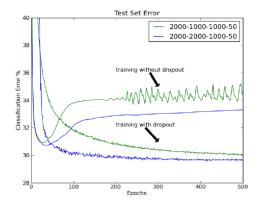
☐ Reuters. Классификация документов.

Входной вектор 2000 элементов.

50 классов.

~200000 тестовых документов.

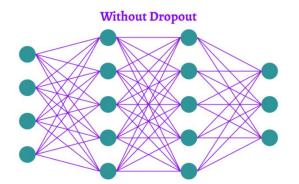


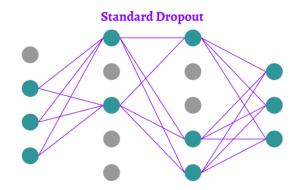


G.E. Hinton et al, Improving neural networks by preventing co-adaptation of feature detectors, 2012

Варианты Dropout. DropConnect.

□ Удаление связей.



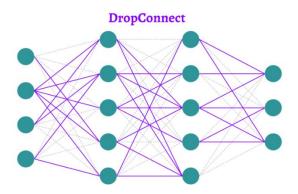


Training Phase:

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

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:

$$\mathbf{y} = (\mathbf{W}\mathbf{x}) \circ \hat{\mathbf{m}}(\mathbf{Z})$$

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

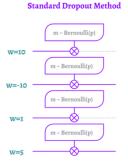
and
$$\boldsymbol{\mu} = p\mathbf{W}\mathbf{x}$$
, $\boldsymbol{\sigma}^2 = p(1-p)(\mathbf{W} \circ \mathbf{W})(\mathbf{x} \circ \mathbf{x})$, $Z \in \mathbb{N}^+$

Варианты Dropout. Standout.

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

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

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



Training Phase:

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

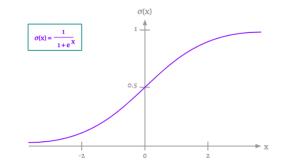
Testing Phase:

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

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







Варианты Dropout. Gaussian Dropout.

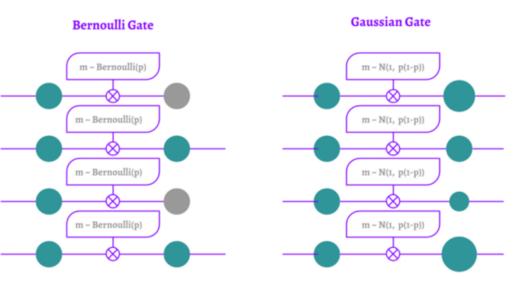
- □ Нейроны не удаляются, а «завешиваются» с помощью нормального распределения.
- □ Выше скорость сходимости.

Training Phase:

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

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:

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







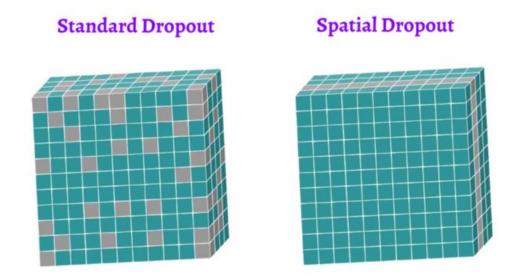






Варианты Dropout. Spatial Dropout.

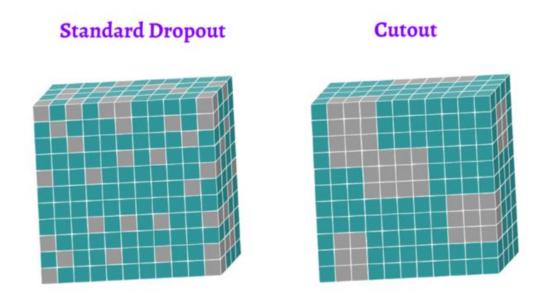
- □ Удаление карт признаков
- ☐ Spatial Dropout применяют для сверточных сетей



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

Dropout B Keras

```
inputs = keras.Input(shape=(180, 180, 3))
x = data augmentation(inputs)
x = layers.Rescaling(1./255)(x)
x = layers.Conv2D(filters=32, kernel size=3, activation="relu")(x)
x = layers.MaxPooling2D(pool size=2)(x)
x = layers.Conv2D(filters=64, kernel size=3, activation="relu")(x)
x = layers.MaxPooling2D(pool size=2)(x)
x = layers.Conv2D(filters=128, kernel size=3, activation="relu")(x)
x = layers.MaxPooling2D(pool size=2)(x)
x = layers.Conv2D(filters=256, kernel size=3, activation="relu")(x)
x = layers.MaxPooling2D(pool size=2)(x)
x = layers.Conv2D(filters=256, kernel size=3, activation="relu")(x)
x = lavers.Flatten()(x)
x = layers.Dropout(0.5)(x)
outputs = layers.Dense(1, activation="sigmoid")(x)
model = keras.Model(inputs=inputs, outputs=outputs)
                                                                ■ Dropout
                                                                ☐ SpatialDropout1D
                                                                ☐ SpatialDropout2D
                                                               ☐ SpatialDropout3D
                                                               ☐ Gaussian Dropout
                                                                ■ AlphaDropout
```

Batch normalization

- lue Нормализация + декорреляция входных данных ightarrow 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). ■ Эффект от Batch normalization:
 - влияние на распределение данных (контроль internal covariance shift)
 - стимуляция градиента.

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 network
   1: N_{\rm BN}^{\rm 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)

 Modify each layer in N<sub>RN</sub> with input x<sup>(k)</sup> to take

                        u^{(k)} instead
   5: end for
   6: Train N_{
m BN}^{
m 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 B, each of
                         size m, and average over them:
                                                                                      E[x] \leftarrow E_{\mathcal{B}}[\mu_{\mathcal{B}}]
                    In N_{BN}^{inf}, replace the transform y = BN_{\gamma,\beta}(x) with
                        y = \frac{\gamma}{\sqrt{\text{Var}[x] + \epsilon}} \cdot x + \left(\beta - \frac{\gamma E[x]}{\sqrt{\text{Var}[x] + \epsilon}}\right)
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}
```

Реализовано через экспоненциальное скользящее среднее, вычисляемое на этапе обучения

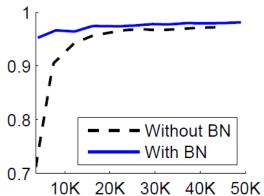
Batch normalization. Примеры.

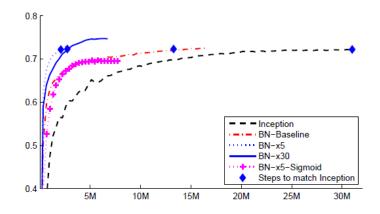
- Классификация MNIST
- □ Сеть 784х100х100х100х10,
- 50000 обучающих примеров

- 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%





Реализация Batch normalization

```
x = layers.Conv2D(32, 3, activation="relu")(x)
x = layers.BatchNormalization()(x)

?
x = layers.Conv2D(32, 3, use_bias=False)(x)
x = layers.BatchNormalization()(x)
x = layers.Activation("relu")(x)
```

Реализация Batch normalization

```
x = layers.Conv2D(32, 3, activation="relu")(x)
x = layers.BatchNormalization()(x)

x = layers.Conv2D(32, 3, use_bias=False)(x)
x = layers.BatchNormalization()(x)
x = layers.Activation("relu")(x)
```

Вопросы

- □ Для чего нужна аугментация данных?
- □ Какие типы Dropout применяют для сверточных сетей?
- □ Почему Batch Normalization повышает эффективность обучения?