

Lecture 7. Convolutional Neural Networks, advanced techniques



Hello, today we continue to talk about Convolutional Neural Networks, advanced optimization methods and architecture.

And again we will start with three questions for five minutes based on the materials of the previous lecture.

└ Five-minutes block

Please, write answers or send photos with them directly to me in private messages here in Teams, so that others cannot read your message. Last time, a lot of people did it, so I believe that this time you all will succeed.

Lecture 7. Convolutional Neural Networks, advanced technics

└ Five-minutes block

Five-minutes block

- Define a convolution operation for neural networks
- Write down the main features of the AlexNet network
- Write the formulas for updating the weights in the Momentum method

Please, write answers or send photos with them directly to me in private messages here in Teams, so that others cannot read your message. Last time, a lot of people did it, so I believe that this time you all will succeed.

Lecture 7. Convolutional Neural Networks, advanced techniques



Let's look at this picture. In many ways, it reflects the entire content of today's lecture.

└ Neural networks optimization methods

A very popular method for improving the convergence of SGD is the momentum accumulation method. In it, from a physical point of view, the derivative of the loss function becomes the acceleration of the change in model parameters, and not the speed as in classical SGD.

Lecture 7. Convolutional Neural Networks, advanced techniques

└ Neural networks optimization methods

Momentum accumulation method
[B.T Polyak, 1964] — exponential moving average of the gradient over $\frac{1}{1-\gamma}$ last iterations:

$$\nu = \gamma\nu + (1 - \gamma)\mathcal{L}'(w)$$

$$w = w - \eta\nu$$



A very popular method for improving the convergence of SGD is the momentum accumulation method. In it, from a physical point of view, the derivative of the loss function becomes the acceleration of the change in model parameters, and not the speed as in classical SGD.

└ AdaGrad

Another interesting approach is the adaptive gradient method. You accumulate squares of gradient vector and then change the model parameters by the component-wise way.

Lecture 7. Convolutional Neural Networks, advanced techniques

└ AdaGrad

$$\begin{aligned}G &= G + \mathcal{L}'(w) \odot \mathcal{L}'(w) \\w &= w - \eta_t \mathcal{L}'(w) \odot \sqrt{G + \varepsilon}\end{aligned}$$

η_t can be fixed, for example $\eta_t = 0.01$.
 \odot, \odot — component-wise multiplication and division of vectors

Another interesting approach is the adaptive gradient method. You accumulate squares of gradient vector and then change the model parameters by the component-wise way.

Lecture 7. Convolutional Neural Networks, advanced techniques

└ AdaGrad

$$\begin{aligned} G &= G + \mathcal{L}_i'(w) \odot \mathcal{L}_i'(w) \\ w &= w - \eta_i \mathcal{L}_i'(w) \odot \sqrt{G + \varepsilon} \end{aligned}$$

η_i can be fixed, for example $\eta_i = 0.01$.

$\odot, \oslash \rightarrow$ component-wise multiplication and division of vectors

Advantages:

- $\mathcal{L}_i'(w)$ — gradient on i -th object
- separate choice of adaptation for each direction
- is suitable for sparse data

Another interesting approach is the adaptive gradient method. You accumulate squares of gradient vector and then change the model parameters by the component-wise way.

Lecture 7. Convolutional Neural Networks, advanced techniques

└ AdaGrad

AdaGrad

$$\begin{aligned} G &= G + \mathcal{L}_i'(w) \odot \mathcal{L}_i'(w) \\ w &= w - \eta_i \mathcal{L}_i'(w) \odot \sqrt{G + \varepsilon} \end{aligned}$$

η_i can be fixed, for example $\eta_i = 0.01$.

$\odot, \oslash \rightarrow$ component-wise multiplication and division of vectors

Advantages:

- $\mathcal{L}_i'(w)$ — gradient on i -th object
- separate choice of adaptation for each direction
- is suitable for sparse data

Disadvantage:

- G is constantly increasing, which can cause learning to stop

Another interesting approach is the adaptive gradient method. You accumulate squares of gradient vector and then change the model parameters by the component-wise way.

Lecture 7. Convolutional Neural Networks, advanced technics

└ RMSProp (running mean square)

The development of AdaGrad is RMSProp — optimization method, which you really have in modern libraries.

Lecture 7. Convolutional Neural Networks, advanced techniques

└ RMSProp (running mean square)

RMSProp (running mean square)

Equalization of rates of change of weights by moving average

$$G = \alpha G + (1 - \alpha) L'_t(w) \odot L'_t(w)$$

$$w = w - \eta_t L'_t(w) \odot \sqrt{G + \epsilon}$$

η_t can be fixed, for example $\eta_t = 0.01$

The development of AdaGrad is RMSProp — optimization method, which you really have in modern libraries.

Lecture 7. Convolutional Neural Networks, advanced techniques

└ RMSProp (running mean square)

RMSProp (running mean square)

Equalization of rates of change of weights by moving average

$$G = \alpha G + (1 - \alpha) L'_t(w) \odot L'_t(w)$$

$$w = w - \eta_t L'_t(w) \odot \sqrt{G + \epsilon}$$

η_t can be fixed, for example $\eta_t = 0.01$

Advantages:

- same as AdaGrad
- and G does not grow uncontrollably

The development of AdaGrad is RMSProp — optimization method, which you really have in modern libraries.

Lecture 7. Convolutional Neural Networks, advanced techniques

└ RMSProp (running mean square)

RMSProp (running mean square)

Equalization of rates of change of weights by moving average

$$G = \alpha G + (1 - \alpha) L'_t(w) \odot L'_t(w)$$

$$w = w - \eta_t L'_t(w) \odot \sqrt{G + \epsilon}$$

η_t can be fixed, for example $\eta_t = 0.01$

Advantages:

- same as AdaGrad
- and G does not grow uncontrollably

Disadvantages:

- at the very beginning of training, G is averaged over zero previous values
- no moment count

The development of AdaGrad is RMSProp — optimization method, which you really have in modern libraries.

Lecture 7. Convolutional Neural Networks, advanced technics

└ Adam (adaptive momentum)

Adam is also one of modern optimization method, which you really may meet in modern applications.

Lecture 7. Convolutional Neural Networks, advanced techniques

└ Adam (adaptive momentum)

$$\begin{aligned}\nu &= \gamma\nu + (1 - \gamma)\mathcal{L}'(\mathbf{w}) & \hat{\nu} &= \nu(1 - \gamma^k)^{-1} \\ G &= \alpha G + (1 - \alpha)\mathcal{L}'(\mathbf{w}) \odot \mathcal{L}'_i(\mathbf{w}) & \hat{G} &= G(1 - \alpha^k)^{-1} \\ w &= w - \eta\hat{\nu} \odot \{\sqrt{\hat{G}} + \varepsilon\}\end{aligned}$$

k — iteration number, $\gamma = 0.9$, $\alpha = 0.999$, $\varepsilon = 10^{-8}$

Adam is also one of modern optimization method, which you really may meet in modern applications.

Lecture 7. Convolutional Neural Networks, advanced techniques

└ Adam (adaptive momentum)

Adam (adaptive momentum)

$$\begin{aligned}\nu &= \gamma\nu + (1 - \gamma)\mathcal{L}'(w) & \hat{\nu} &= \nu(1 - \gamma^k)^{-1} \\ G &= \alpha G + (1 - \alpha)\mathcal{L}'(w) \odot \mathcal{L}'_i(w) & \hat{G} &= G(1 - \alpha^k)^{-1} \\ w &= w - \eta\hat{\nu} \odot (\sqrt{\hat{G}} + \varepsilon)\end{aligned}$$

k — iteration number, $\gamma = 0.9$, $\alpha = 0.999$, $\varepsilon = 10^{-8}$

Advantages:

- are almost the same as RMSProp
- calibration $\hat{\nu}, \hat{G}$ increases ν, G on the first iterations
- is count of moment

Adam is also one of modern optimization method, which you really may meet in modern applications.

Lecture 7. Convolutional Neural Networks, advanced techniques

└ Nadam (Nesterov-accelerated adaptive momentum)

Same formulas for v, β, G, \hat{G}
 $w = w - \eta_t \left(\gamma \hat{p} + \frac{1-\gamma}{1-\gamma^t} \hat{G}(w) \right) \odot (\sqrt{G + \varepsilon})$
 k — iteration number, $\gamma = 0.9$, $\alpha = 0.999$, $\varepsilon = 10^{-8}$

T. Dozat: Incorporating Nesterov Momentum into Adam. ICLR-2016.

And Nadam (Nesterov-accelerated adaptive momentum) — is modern modification of Adam method.

Lecture 7. Convolutional Neural Networks, advanced technics

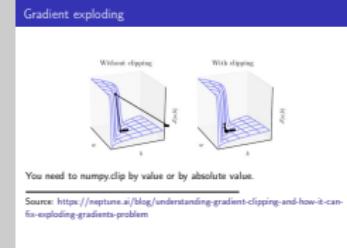
└ Convergence example

Source: <http://www.denizyuret.com/2015/03/alec-radford-s-animations-for.html>

Let's look on example of convergence of considered optimization methods.

Lecture 7. Convolutional Neural Networks, advanced techniques

└ Gradient exploding



There is typical problem in neural network optimization, it is gradient exploding.

Lecture 7. Convolutional Neural Networks, advanced technics

└ Regularization Methods

Dropout is a method of random shutdowns of neurons. It has two stages: learning stage and application stage.

Lecture 7. Convolutional Neural Networks, advanced technics

└ Regularization Methods

Regularization Methods

Dropout is a method of random shutdowns of neurons

Dropout is a method of random shutdowns of neurons. It has two stages:
learning stage and application stage.

Lecture 7. Convolutional Neural Networks, advanced techniques

└ Regularization Methods

Regularization Methods

Dropout is a method of random shutdowns of neurons
Learning stage: make the gradient step $\mathcal{L}_t(w) \rightarrow \min_w$, disable the b -th neuron of the ℓ -th layer with probability p_t :
 $x_{tb}^{\ell+1} = \tilde{x}_{tb}^{\ell} \sigma_{tb}(\sum_j w_{jb} x_{j b}^{\ell}), \quad P(x_{tb}^{\ell} = 0) = p_t$

Lecture 7. Convolutional Neural Networks, advanced techniques

└ Regularization Methods

Regularization Methods

Dropout is a method of random shutdowns of neurons

Learning stage: make the gradient step $\mathcal{L}_t(w) \rightarrow \min_w$ disable the b -th neuron of the t -th layer with probability p_t :
 $x_{jb}^{t+1} = \sum_j w_{jt} x_{jt}^t, \quad P(x_{jb}^t = 0) = p_t$

Application stage: turn on all neurons, but with a correction:
 $x_{jb}^{-1} = (1 - p_t) x_{jb} + \sum_j w_{jt} x_{jt}^t$

Lecture 7. Convolutional Neural Networks, advanced techniques

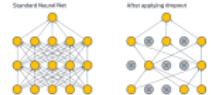
└ Regularization Methods

Regularization Methods

Dropout is a method of random shutdowns of neurons

Learning stage: make the gradient step $\mathcal{L}_t(w) \rightarrow \min_w$ disable the h -th neuron of the t -th layer with probability p_t :
 $x_{jh}^{t+1} = \sum_j w_{jh} x_{jh}^t, \quad P(x_{jh}^t = 0) = p_t$

Application stage: turn on all neurons, but with a correction:
 $x_{jh}^{-1} = (1 - p_t) x_{jh}^t \sum_j w_{jh} x_{jh}^t$



Dropout is a method of random shutdowns of neurons. It has two stages: learning stage and application stage.

Lecture 7. Convolutional Neural Networks, advanced techniques

└ Batch normalization

Batch normalization

$B = \{x_i\}$ — mini-batch of whole dataset

Averaging the $\mathcal{L}(w)$ gradients over the batch speeds up the convergence.

$B^l = \{a_i^l\}$ — vectors of x_i objects at the output of the l -th layer

Let's go on to the next important regularization method. It is Batch normalization method.

Lecture 7. Convolutional Neural Networks, advanced techniques

└ Batch normalization

Batch normalization

$B = \{x_i\}$ — mini-batch of whole dataset

Averaging the $\mathcal{L}(w)$ gradients over the batch speeds up the convergence.

$B^l = \{u_i^l\}$ — vectors of x_i objects at the output of the l -th layer

1. Normalize each j -th component of the vector u_i^l by the batch

$$u_i^l = \frac{u_i^l - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}, \quad \mu_j = \frac{1}{|B|} \sum_{k \in B} u_k^l; \quad \sigma_j^2 = \frac{1}{|B|} \sum_{k \in B} (u_k^l - \mu_j)^2$$

Let's go on to the next important regularization method. It is Batch normalization method.

Lecture 7. Convolutional Neural Networks, advanced techniques

└ Batch normalization

Batch normalization

$B = \{x_i\}$ — mini-batch of whole dataset

Averaging the $\mathcal{L}(w)$ gradients over the batch speeds up the convergence.

$B^l = \{u_i^l\}$ — vectors of x_i objects at the output of the l -th layer

1. Normalize each j -th component of the vector u_i^l by the batch

$$\hat{u}_{ij}^l = \frac{u_{ij}^l - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}, \quad \mu_j = \frac{1}{|B|} \sum_{i \in B} u_{ij}^l; \quad \sigma_j^2 = \frac{1}{|B|} \sum_{i \in B} (u_{ij}^l - \mu_j)^2$$

2. Add a linear layer with custom weights

$$\tilde{u}_{ij}^l = \gamma_j^l \hat{u}_{ij}^l + \beta_j^l$$

Let's go on to the next important regularization method. It is Batch normalization method.

Lecture 7. Convolutional Neural Networks, advanced techniques

└ Batch normalization

Let's go on to the next important regularization method. It is Batch normalization method.

Batch normalization

$B = \{x_i\}$ — mini-batch of whole dataset

Averaging the $\mathcal{L}(w)$ gradients over the batch speeds up the convergence.

$B^l = \{u_i^l\}$ — vectors of x_i objects at the output of the l -th layer

1. Normalize each j -th component of the vector u_i^l by the batch

$$\hat{u}_{ij}^l = \frac{u_{ij}^l - \mu_j}{\sqrt{\sigma_j^2 + \epsilon}}, \quad \mu_j = \frac{1}{|B|} \sum_{i \in B} u_{ij}^l; \quad \sigma_j^2 = \frac{1}{|B|} \sum_{i \in B} (u_{ij}^l - \mu_j)^2$$

2. Add a linear layer with custom weights:

$$\tilde{u}_{ij}^l = \gamma_j^l \hat{u}_{ij}^l + \beta_j^l$$

3. Parameters γ_j^l, β_j^l are configured by BackProp

Lecture 7. Convolutional Neural Networks, advanced techniques

└ Initial Approximation of the Weights



Also of course very important to choose consistently the Initial Approximation of the Weights.

Lecture 7. Convolutional Neural Networks, advanced techniques

└ Monitoring changes in variances

Monitoring changes in variances
Both of neuron values and backpropagation gradients

Equalization of variances of both neurons and gradients in different layers
— Xavier initialization for tasks:
$$w_j \sim U\left[-\frac{6}{\sqrt{n_{in}} + \sqrt{n_{out}}}, \frac{6}{\sqrt{n_{in}} + \sqrt{n_{out}}}\right]$$

 n_{in}, n_{out} — the number of neurons in the previous and current layers, respectively
The proof and details are in the SHAD textbook (in Russian yet).

In fact, the consistently rule is quite simple — you need to make sure that the variance of the layers is preserved both during the forward and backward pass through the neural network.

Lecture 7. Convolutional Neural Networks, advanced techniques

└ Monitoring changes in variances

Monitoring changes in variances
Both of neuron values and backpropagation gradients

Equalization of variances of both neurons and gradients in different layers
— Xavier initialization for tasks:

$$w_j \sim U \left[-\frac{6}{\sqrt{n_{in}} + \sqrt{n_{out}}} \frac{6}{\sqrt{n_{in}} + \sqrt{n_{out}}} \right]$$

n_{in}, n_{out} — the number of neurons in the previous and current layers, respectively

The proof and details are in the SHAD textbook (in Russian yet).

Question

What to do in case of ReLU activation?

In fact, the consistently rule is quite simple — you need to make sure that the variance of the layers is preserved both during the forward and backward pass through the neural network.

Lecture 7. Convolutional Neural Networks, advanced techniques

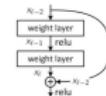
└ ResNet — Residual Net

ResNet — Residual Net

Skip-connection of the ℓ layer with the previous $\ell - d$ layer:

$$x_\ell = \sigma(W_\ell x_{\ell-1}) + x_{\ell-d}$$

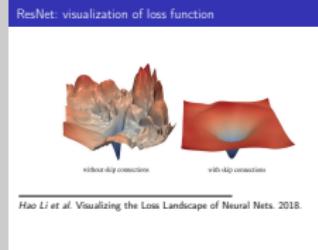
The ℓ layer learns not the new vector representation x_ℓ , but its increment $x_\ell - x_{\ell-d}$



Great, now let's talk a bit about modern neural network architectures.

Lecture 7. Convolutional Neural Networks, advanced techniques

└ ResNet: visualization of loss function



This is famous image from a paper that almost became a meme in deep learning :)

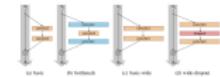
Lecture 7. Convolutional Neural Networks, advanced techniques

└ WideResNet — WRN

WideResNet — WRN

Disadvantages of ResNet

- is too deep — up to 1000 layers :)
- take a long time to train to SoTA



- a simple WRN network of 16 layers defeated all ResNets
- WRN-40-4 (i.e. 40 layers and 4x width) trains 8 times faster than ResNet-1001

S. Zagoruyko, N. Komodakis Wide Residual Networks, 2017

And the final architecture for today's lecture WideResNet, which I recommend paying attention to when solving the competition from the second homework. It has already been posted on github and in teams.

Temporary page!

\LaTeX was unable to guess the total number of pages correctly. A
was some unprocessed data that should have been added to the
this extra page has been added to receive it.

If you rerun the document (without altering it) this surplus page
away, because \LaTeX now knows how many pages to expect for the
document.