

Training DNN

HOML – chapter11

TAVE Research DL001

Changdae Oh

2021. 01. 10

Contents

The challenges of training Deep Neural Networks

- Vanishing gradient / Exploding gradient
- Insufficient data, Difficulties in labeling
- Training can be extremely slow
- Risk of overfitting

1. Gradient vanishing/exploding
2. Reusing pre-trained layer
3. Fast optimizer
4. Avoid overfitting by regularization
5. Summary



Contents

1. Gradient vanishing/exploding
2. Reusing pre-trained layer
3. Fast optimizer
4. Avoid overfitting by regularization
5. Summary

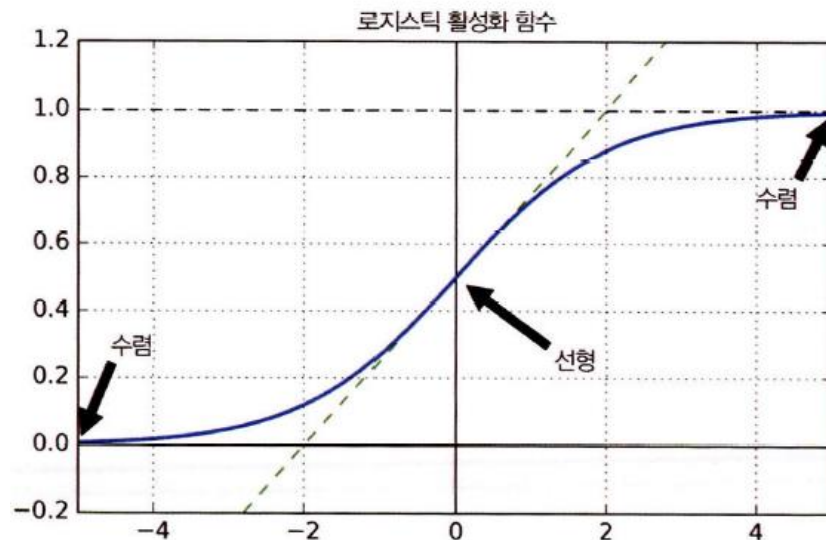
Gradient vanishing/exploding

Vanishing gradient

- When training DNN, the gradient gradually decreases as it gets closer to the input layer in the backpropagation process

Exploding gradient

- On the contrary, as it approaches the input layer, the gradient gradually increases and diverges.



How can we solve this problem?

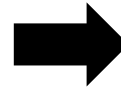
- 1) Initializing strategy
- 2) Activation functions
- 3) Batch Normalization
- 4) Gradient Clipping

Gradient vanishing/exploding

1. Initializing strategy

Past weight initialization strategy

: Sampling from
Normal distribution (Mean : 0, Var : 1)



Glorot initialization (Xavier initialization)

Normal distribution (Mean : 0, Var : $\frac{1}{fan_{avg}}$)

uniform distribution ($-\sqrt{\frac{3}{fan_{avg}}}$, $+\sqrt{\frac{3}{fan_{avg}}}$)

LuCun initialization

Just replace fan-avg with fan-in

He initialization

Just double scale variance

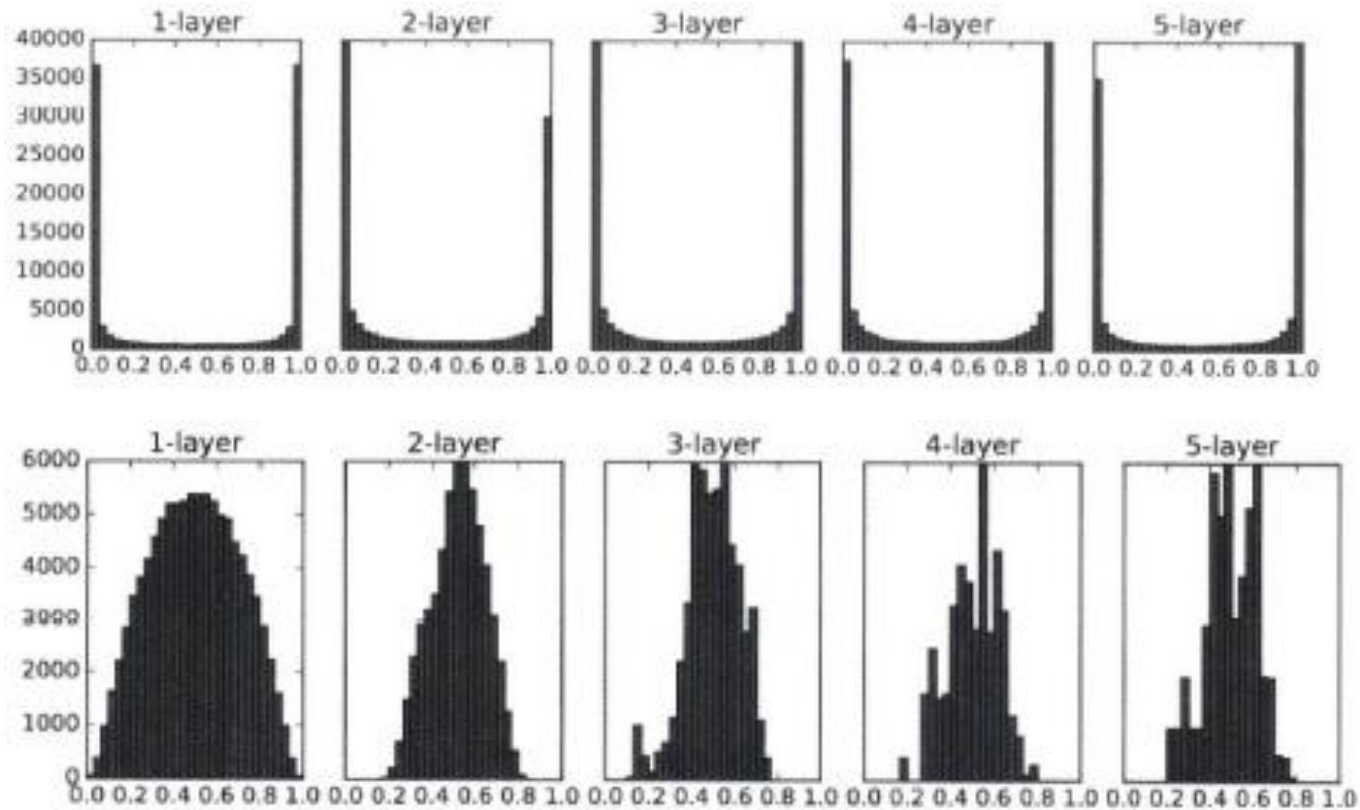
strategy	Activation	Variance(Normal dist.)
Glorot	sigmoid, tanh, softmax	$1 / fan_{avg}$
He	ReLU s	$2 / fan_{in}$
LuCun	SELU	$1 / fan_{avg}$

fan-in : # of input units to weight matrix

fan-out : # of output units from weight matrix

Gradient vanishing/exploding

1. Initializing strategy



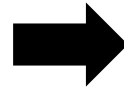
Ex) sigmoid & Glorot init

Gradient vanishing/exploding

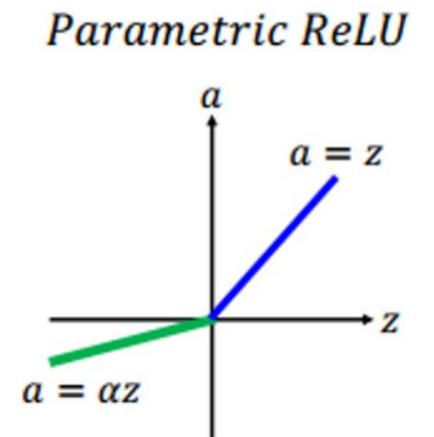
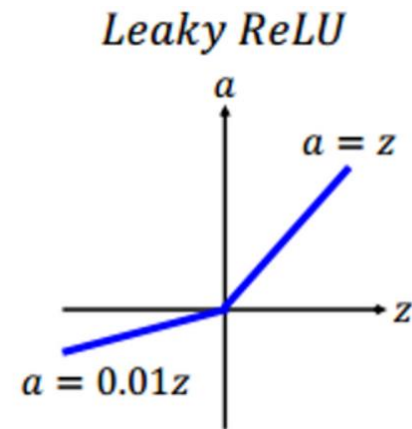
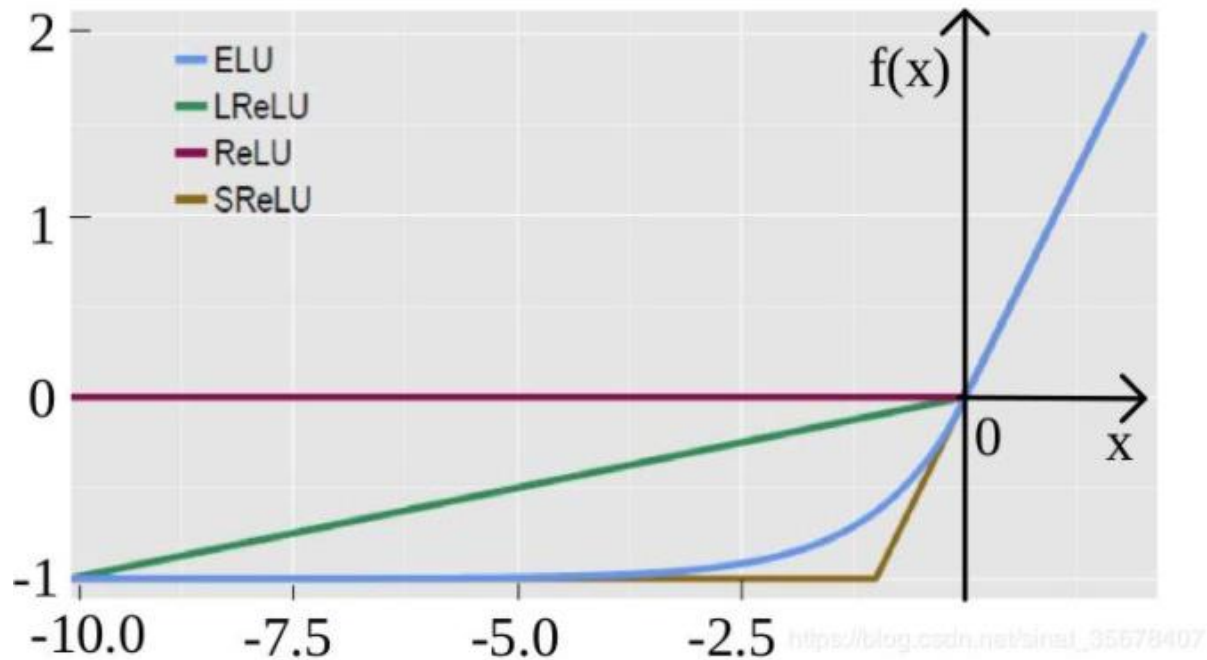
2. non-convergent activation function

< ReLU >

- Not converged to specific positive value
- dying ReLU



- LeakyReLU / RReLU / PReLU
- ELU / SELU

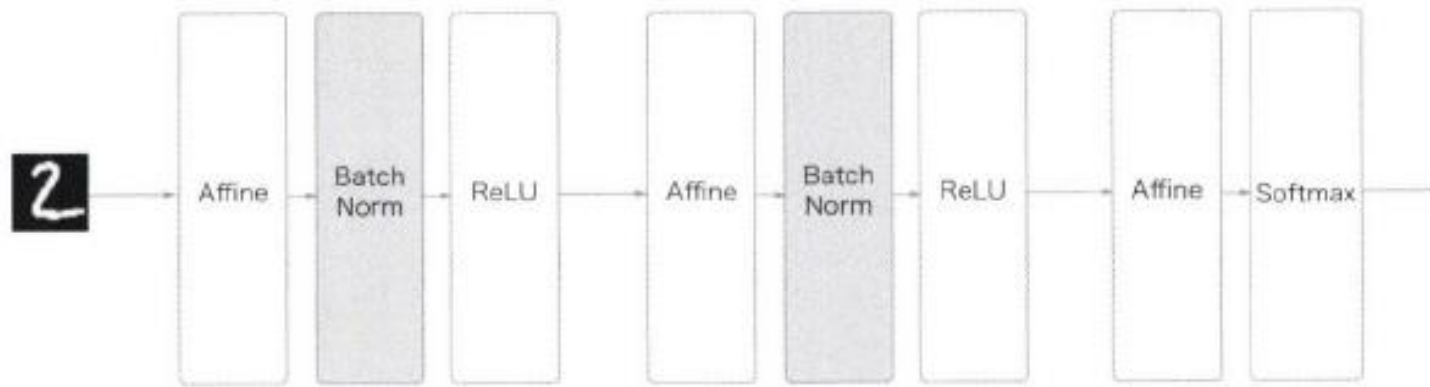


Gradient vanishing/exploding

3. Batch Normalization

- Align the input to the origin and normalize it
- In each layer, two new parameters are used to scale and shift the results.

➡ *Force each layer to spread the activations moderately*



advantage

- Learning faster
- Rubust to initializing
- Restraint overfitting

disadvantage

- Increase complexity
- Computing resource

algorithm

1. $\mu_B = \frac{1}{m_B} \sum_{i=1}^{m_B} \mathbf{x}^{(i)}$
2. $\sigma_B^2 = \frac{1}{m_B} \sum_{i=1}^{m_B} (\mathbf{x}^{(i)} - \mu_B)^2$
3. $\hat{\mathbf{x}}^{(i)} = \frac{\mathbf{x}^{(i)} - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$
4. $\mathbf{z}^{(i)} = \gamma \otimes \hat{\mathbf{x}}^{(i)} + \beta$

Gradient vanishing/exploding

4. Gradient Clipping

- If the norm of the gradient is greater than some threshold, scale it down before applying SGD update

➡ *Take a step in the same direction, but a smaller step*

Algorithm 1 Pseudo-code for norm clipping

```
 $\hat{\mathbf{g}} \leftarrow \frac{\partial \mathcal{E}}{\partial \theta}$   
if  $\|\hat{\mathbf{g}}\| \geq threshold$  then  
     $\hat{\mathbf{g}} \leftarrow \frac{threshold}{\|\hat{\mathbf{g}}\|} \hat{\mathbf{g}}$   
end if
```



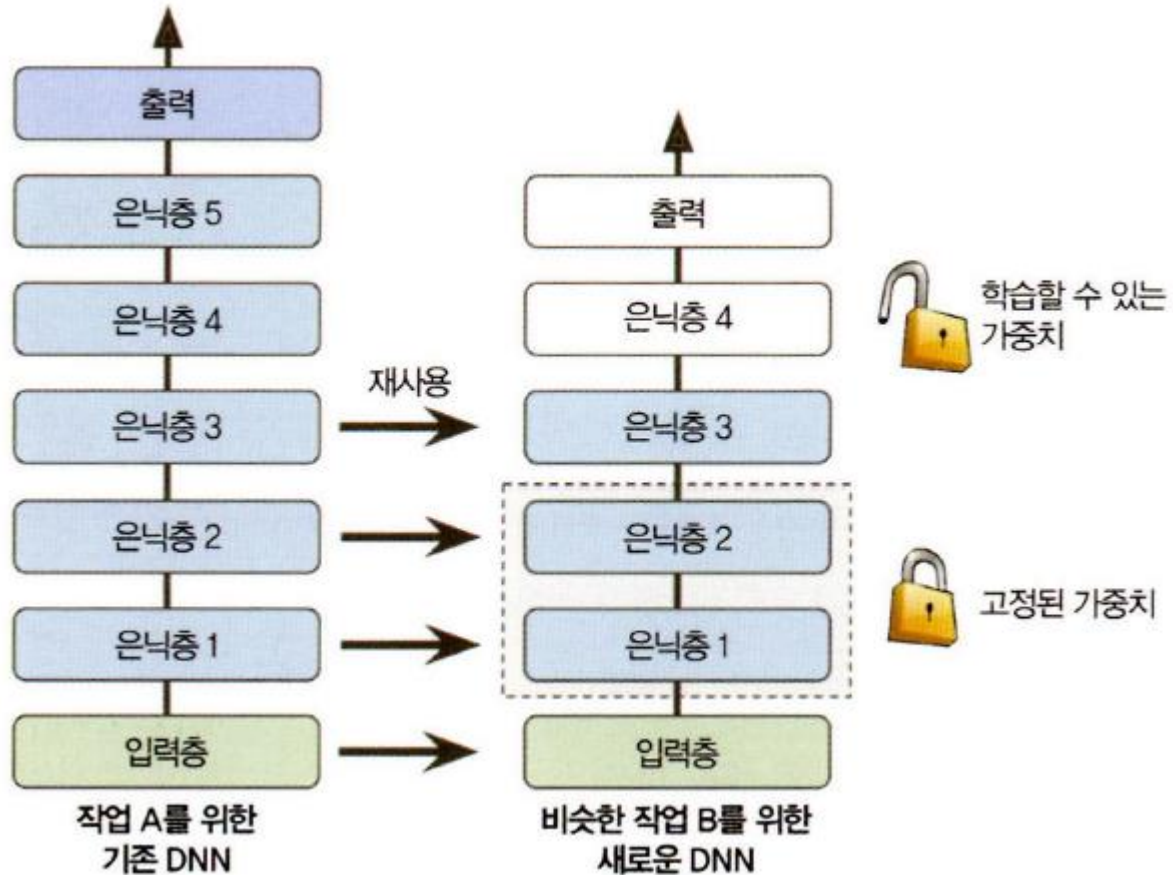
Contents

1. Gradient vanishing/exploding
2. Reusing pre-trained layer
3. Fast optimizer
4. Avoid overfitting by regularization
5. Summary

Reusing pre-trained layer

Transfer learning

- Find out if there is a system that handles a similar type of problem to be solved
- Reuse the lower layers of the neural network.



- ✓ Learning faster
- ✓ Significantly reduce the amount of train data required

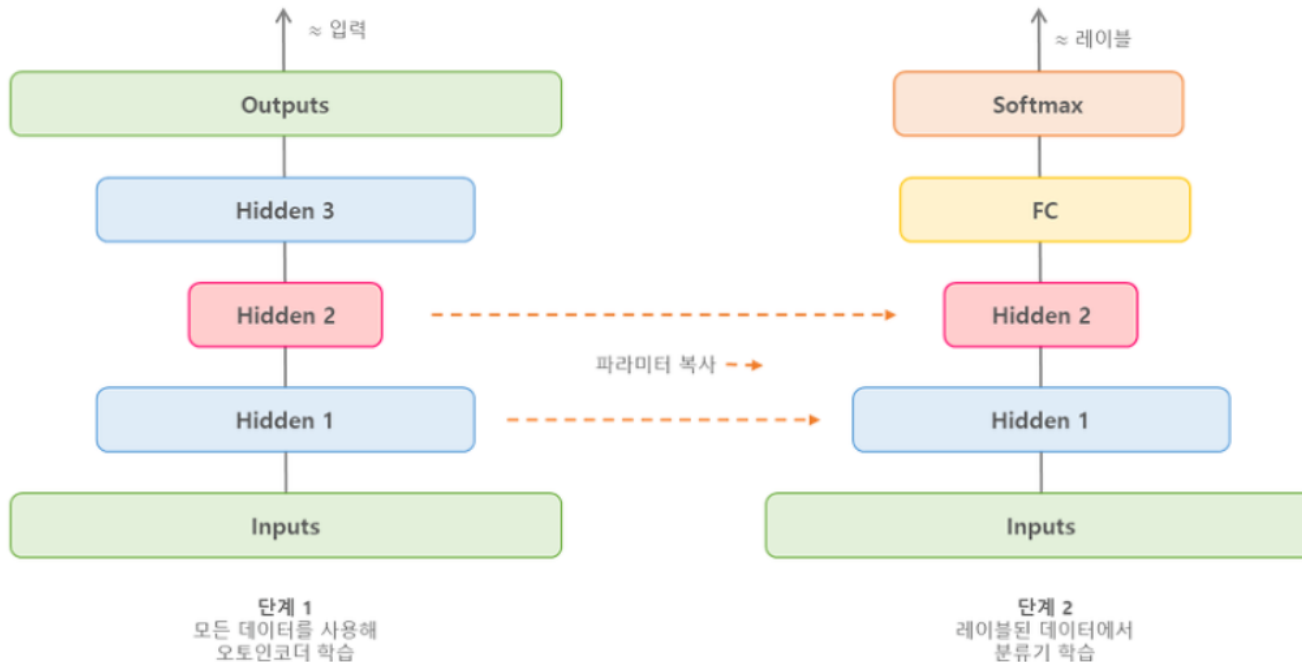
< Process >

1. Freeze all reusable layers.
2. Fitting the model, evaluate performance
3. Compare performance by melting the freezing of one or two layers at the top of the frozen layers and updating the weights via backpropagation

Reusing pre-trained layer

Unsupervised pretraining

- don't have enough labeled training data?



< Process >

1. Train an unsupervised learning model on unlabeled data. (auto encoder, GAN)
2. Reuse lower layers of those models, On top of that, Add hidden layers and an output layer suitable for the new task.
3. Fine-tuning the final network with supervised learning using labeled training samples

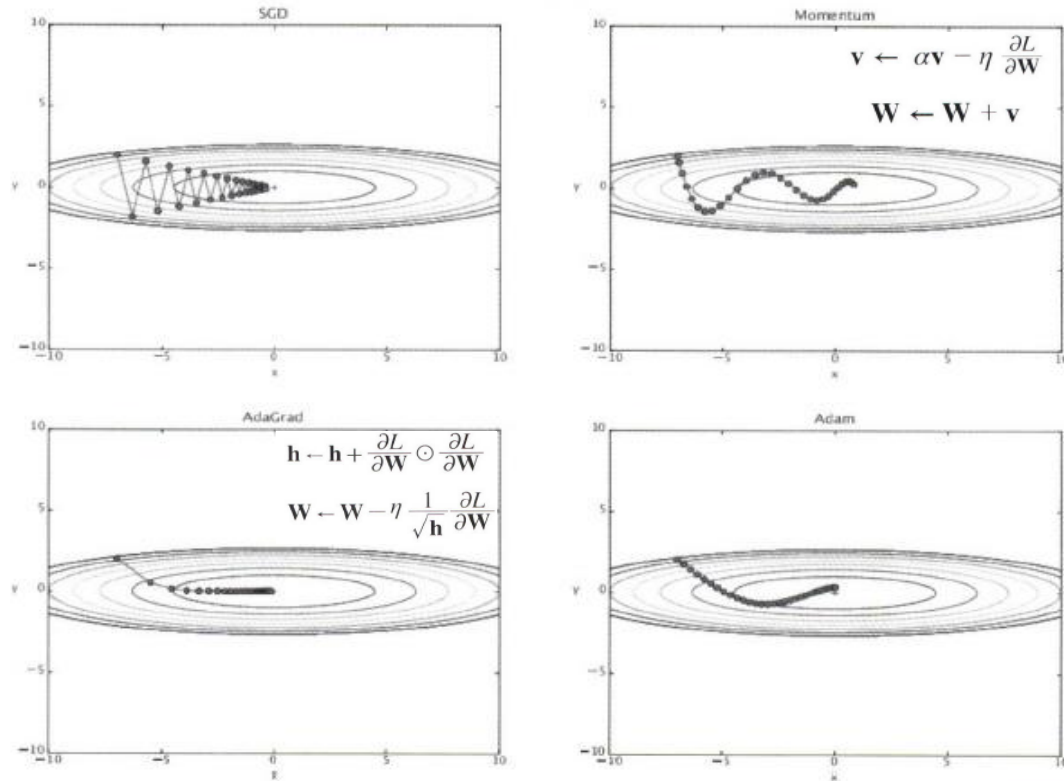
or Pretrain in sub-task !!



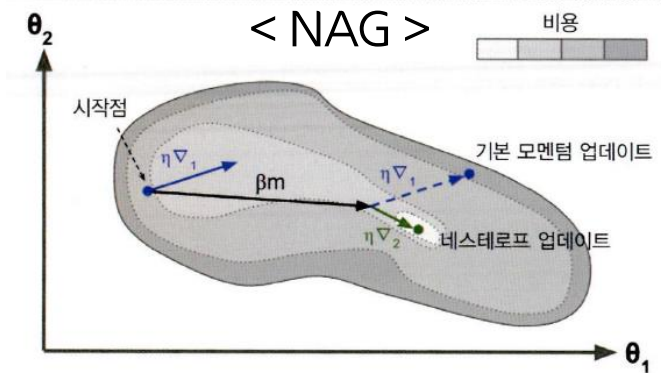
Contents

1. Gradient vanishing/exploding
2. Reusing pre-trained layer
- 3. Fast optimizer**
4. Avoid overfitting by regularization
5. Summary

Fast optimizer



- Momentum / Nesterov accelerated gradient
- Adagrad / RMSProp
- Adam / AdaMax / Nadam

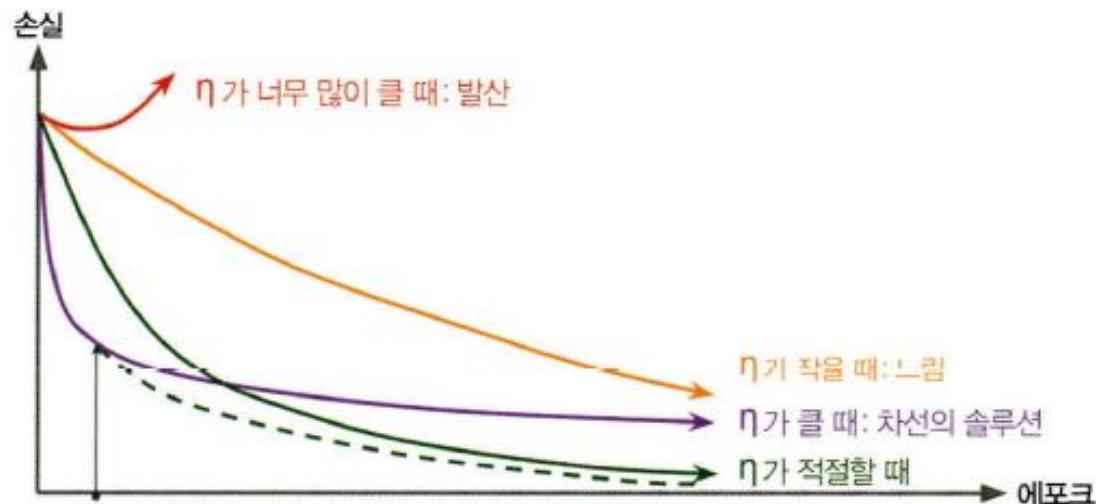


Optimizer	Convergence Speed	Convergence Quality
SGD	*	***
Momentum	**	***
NAG	**	***
Adagrad	***	* (stop too early)
RMSprop	***	** or ***
Adam	***	** or ***
Nadam	***	** or ***
AdaMax	***	** or ***

Fast optimizer

Learning rate Scheduling – LR decay method

Starting with a large learning rate
and decaying the learning rate when the learning rate is slow,
the solution can be found faster than the optimal fixed learning rate



- Power scheduling
- Exponential scheduling
- Piecewise constant scheduling
- Performance scheduling
- 1 cycle scheduling



Contents

1. Gradient vanishing/exploding
2. Reusing pre-trained layer
3. Fast optimizer
- 4. Avoid overfitting by regularization**
5. Summary

Avoid overfitting by regularization

L1, L2 Regularization

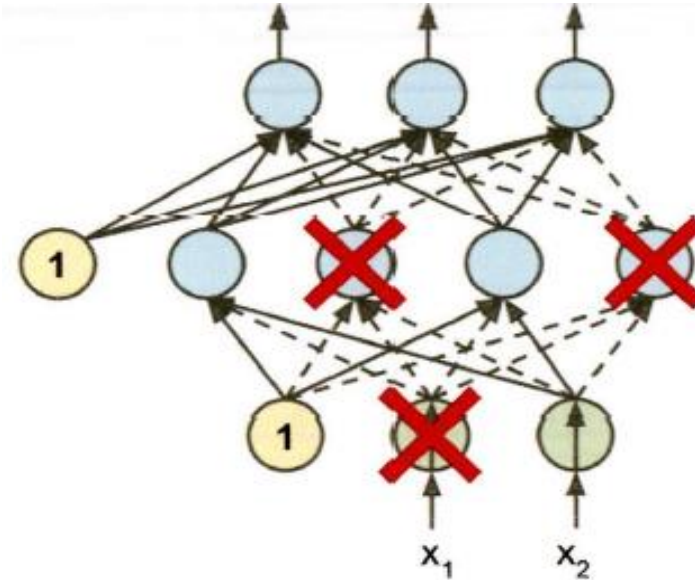
Norm term $\|\mathbf{x}\|_p := \left(\sum_{i=1}^n |x_i|^p \right)^{1/p}$

L1
$$Cost = \frac{1}{n} \sum_{i=1}^n \{L(y_i, \hat{y}_i) + \frac{\lambda}{2} |w|\}$$

L2
$$Cost = \frac{1}{n} \sum_{i=1}^n \{L(y_i, \hat{y}_i) + \frac{\lambda}{2} |w|^2\}$$

Dropout

With hyperparameter p



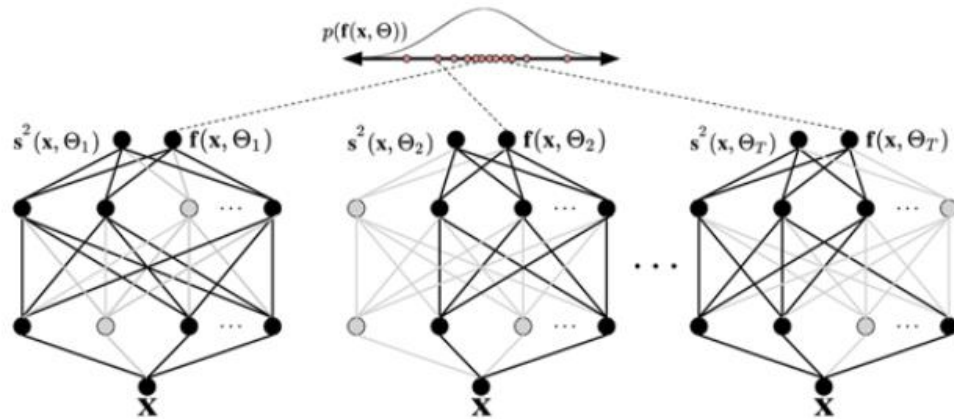
- Make a model **more stable**
- Can be seen as a kind of **ensemble** of individual NNs

End of training, the remaining connection weights are **multiplied by the keep probability**.

There is no action during the testing.

Avoid overfitting by regularization

Monte Carlo Dropout



Typical dropout techniques are applied only during training, and predictions are made using all neurons in the model during testing.



What if dropout is applied to both training and testing processes?

Max-norm regularization

- For each neuron, norm of the input connection weight is limited

$$\|\mathbf{w}\|_2 \leq r$$

- At the end of each training step, the norm of weights is calculated and scaled.

$$\mathbf{w} \leftarrow \mathbf{w} \frac{r}{\|\mathbf{w}\|_2}$$

Summary

hyper-parameter	nice default value
initializing	He
activation	ELU
normalizing	if network is deep - BN
regularization	early stopping / L2
optimizer	momentum(or RMSProp, Nadam)
LR schedule	1-cycle

Exception

- If you need a 'sparse' model, use L1 regularization.
- If a model with a fast response is needed, reduce the number of layers and merge the batch normalization layer into the previous layer.
or use LeakyReLU / ReLU
- Use MCdropout for applications that are risk-sensitive, and whose speed is not very important.

DL from Scratch

Step 8 ~ 16

- Big Picture
- [goal 1] Automatic Differentiation
- [goal 2] Code Improvement

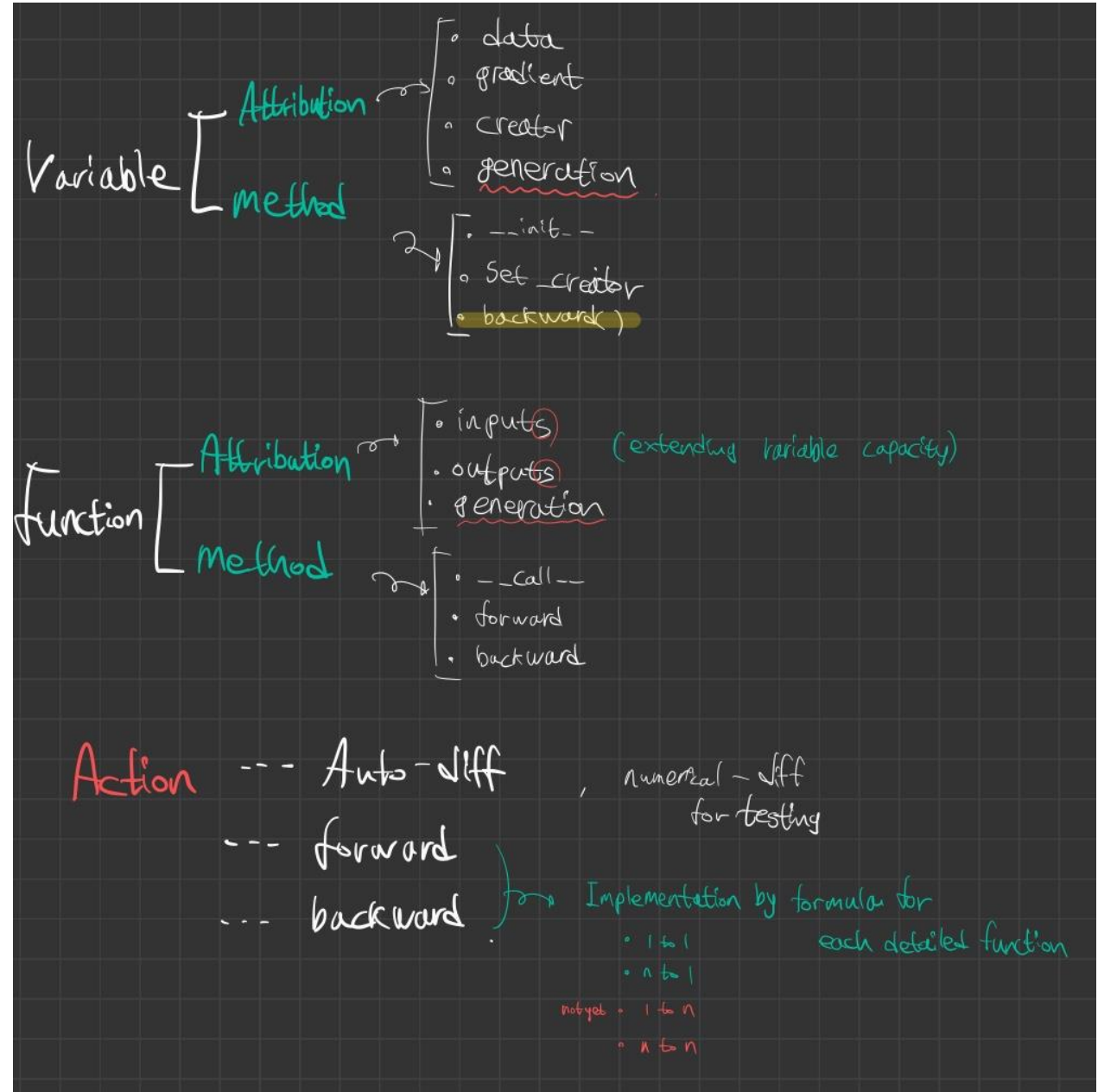
Review & Big Picture

So far. (~ step 16)

What subcomponents are needed
to implement neural networks?

Variable

Function



Automatic Differentiation



< Logic >

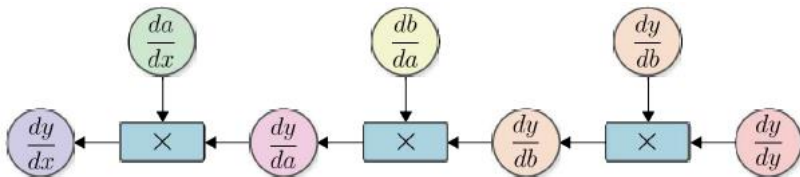
$$\frac{dy}{dx} = \left(\left(\left(\frac{dy}{dy} \frac{dy}{db} \right) \frac{db}{da} \right) \frac{da}{dx} \right)$$

① $\frac{dy}{dy} \frac{dy}{db} = \frac{dy}{db}$

② $\frac{dy}{db} \frac{db}{da} = \frac{dy}{da}$

③ $\frac{dy}{da} \frac{da}{dx} = \frac{dy}{dx}$

< Computation graph >



```
4 class Variable:
5     def __init__(self, data):
6         if data is not None:
7             if not isinstance(data, np.ndarray):
8                 raise TypeError('{} is not supported'.format(type(data)))
9
10        self.data = data
11        self.grad = None
12        self.creator = None
13
14    def set_creator(self, func):
15        self.creator = func
16
17    def backward(self):
18        if self.grad is None:
19            self.grad = np.ones_like(self.data)
20
21        funcs = [self.creator]
22        while funcs:
23            f = funcs.pop()
24            x, y = f.input, f.output
25            x.grad = f.backward(y.grad)
26
27            if x.creator is not None:
28                funcs.append(x.creator)
```

```
31 def as_array(x):
32     if np.isscalar(x):
33         return np.array(x)
34     return x
35
36 class Function:
37     def __call__(self, input):
38         x = input.data
39         y = self.forward(x)
40         output = Variable(as_array(y))
41         output.set_creator(self)
42         self.input = input
43         self.output = output
44         return output
45
46 def forward(self, x):
47     raise NotImplementedError()
48
49 def backward(self, gy):
50     raise NotImplementedError()
```

```
53 class Square(Function):
54     def forward(self, x):
55         y = x ** 2
56         return y
57
58     def backward(self, gy):
59         x = self.input.data
60         gx = 2 * x * gy
61         return gx
```

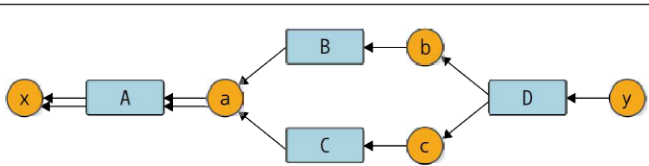
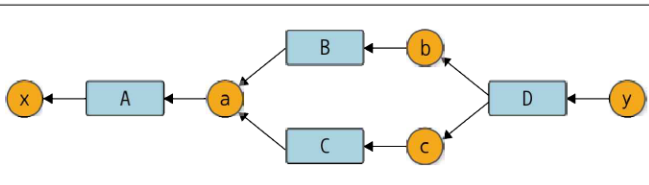
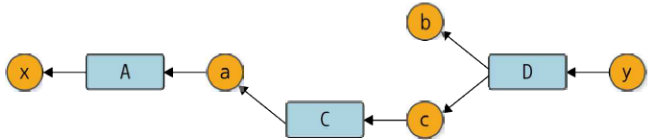
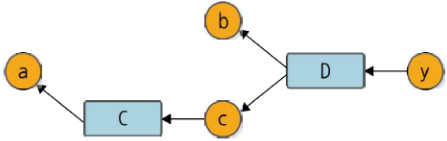
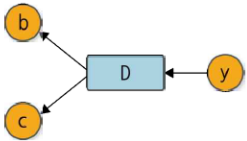
```
67 class Exp(Function):
68     def forward(self, x):
69         y = np.exp(x)
70         return y
71
72     def backward(self, gy):
73         x = self.input.data
74         gx = np.exp(x) * gy
75         return gx
```

Code Improvement

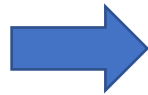
< problems >

1. Single input & output
2. Reuse same variable
3. Order of calculation

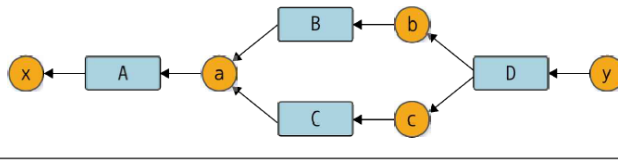
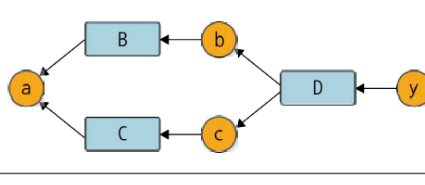
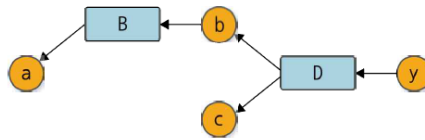
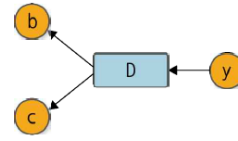
wrong



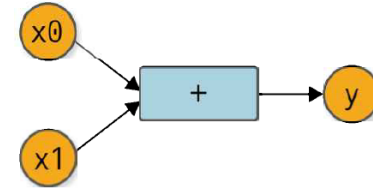
(3)



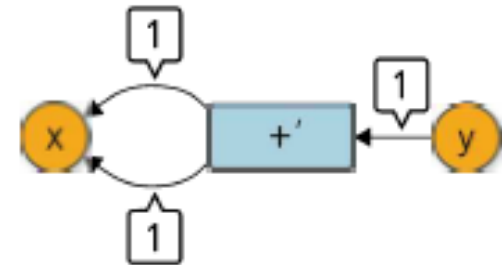
right



(1)



(2)



Code Improvement

```

63 class Function:
64     def __call__(self, *inputs):
65         xs = [x.data for x in inputs] # multiple inputs
66         ys = self.forward(*xs)
67         if not isinstance(ys, tuple): # Tupleizing outputs
68             ys = (ys,)
69         outputs = [Variable(as_array(y)) for y in ys]
70
71         self.generation = max([x.generation for x in inputs]) # setting 'generation' of function
72         for output in outputs:
73             output.set_creator(self)
74         self.inputs = inputs
75         self.outputs = outputs
76         return outputs if len(outputs) > 1 else outputs[0] # to return 1 or many outputs correctly
77
78     def forward(self, xs):
79         raise NotImplementedError()
80
81     def backward(self, gys):
82         raise NotImplementedError()
83
84
85 class Add(Function):
86     # multiple inputs
87     def forward(self, x0, x1):
88         y = x0 + x1
89         return y
90
91     # multiple backward gradients
92     def backward(self, gy):
93         return gy, gy
94
95 def add(x0, x1):
96     return Add()(x0, x1)
97
98 class Square(Function):
99     def forward(self, x):
100         y = x ** 2
101         return y
102
103     def backward(self, gy):
104         x = self.inputs[0].data # to deal with tupleized single input
105         gx = 2 * x * gy
106         return gx
107
108

```

*Tupleizing
All single elements*

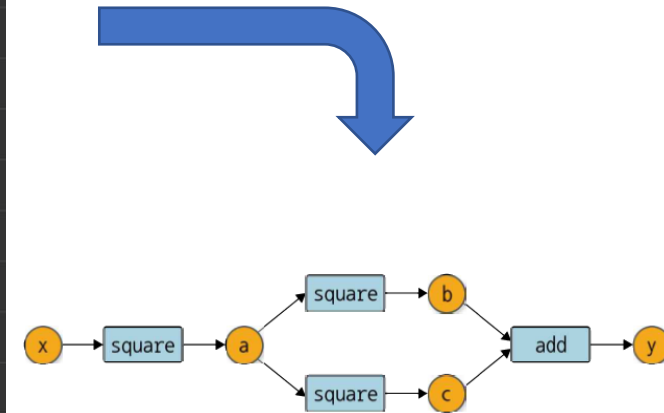
*forward
2 → 1
backward
2 ← 1*

```

4 class Variable:
5     def __init__(self, data):
6         if data is not None:
7             if not isinstance(data, np.ndarray):
8                 raise TypeError('{} is not supported'.format(type(data)))
9
10        self.data = data
11        self.grad = None
12        self.creator = None
13        # use 'generation' variable to control the flow of backward gradients
14        self.generation = 0
15
16    def set_creator(self, func):
17        self.creator = func
18        # variable & function both have 'generation'
19        self.generation = func.generation + 1
20
21    def backward(self):
22        if self.grad is None:
23            self.grad = np.ones_like(self.data)
24
25        # funcs & seen_set : space that functions will be stacked
26        funcs = []
27        seen_set = set() # set data type -> with no duplicates elements
28
29    def add_func(f):
30        if f not in seen_set:
31            funcs.append(f)
32            seen_set.add(f)
33            funcs.sort(key = lambda x: x.generation) # sorting by generation
34        add_func(self.creator)
35
36    while funcs:
37        f = funcs.pop() # taken out sequentially from the higher generation
38        gys = [output.grad for output in f.outputs] # multiple gradients
39        gxs = f.backward(*gys)
40        if not isinstance(gxs, tuple):
41            gxs = (gx,)
42
43        for x, gx in zip(f.inputs, gxs):
44            # prevent overwriting of gradient when using the same variable instance
45            if x.grad is None:
46                x.grad = gx
47            else:
48                x.grad = x.grad + gx
49
50            if x.creator is not None:
51                add_func(x.creator)
52
53    # prevent accumulating of gradient when using the same variable instance
54    def cleargrad(self):
55        self.grad = None

```

*ex) $y = \text{add}(x, x)$
 $= 2x$
 $\frac{dy}{dx} = 2$*



$$y = (x^2)^2 + (x^2)^2 = 2x^4$$

```

x = Variable(np.array(2.0))
a = square(x)
y = add(square(a), square(a))
y.backward()

print(y.data)
print(x.grad)

```

```

32.0
64.0

```