# CS231n Lecture 7

- BOAZ 10기 박성현
- BOAZ 11기 김태희
- BOAZ 11기 홍지민

- Parameter update schemes
- Learning rate schedules
- Gradient checking
- Regularization (Dropout etc.)
- Evaluation (Ensembles etc.)
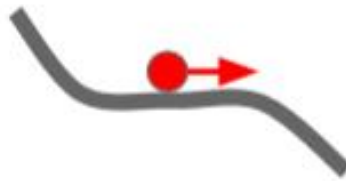- Transfer learning / fine-tuning
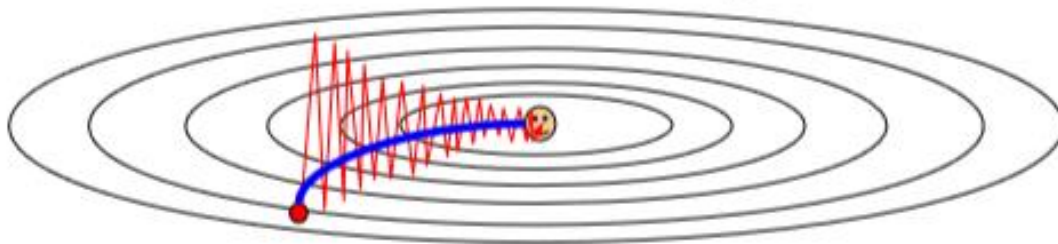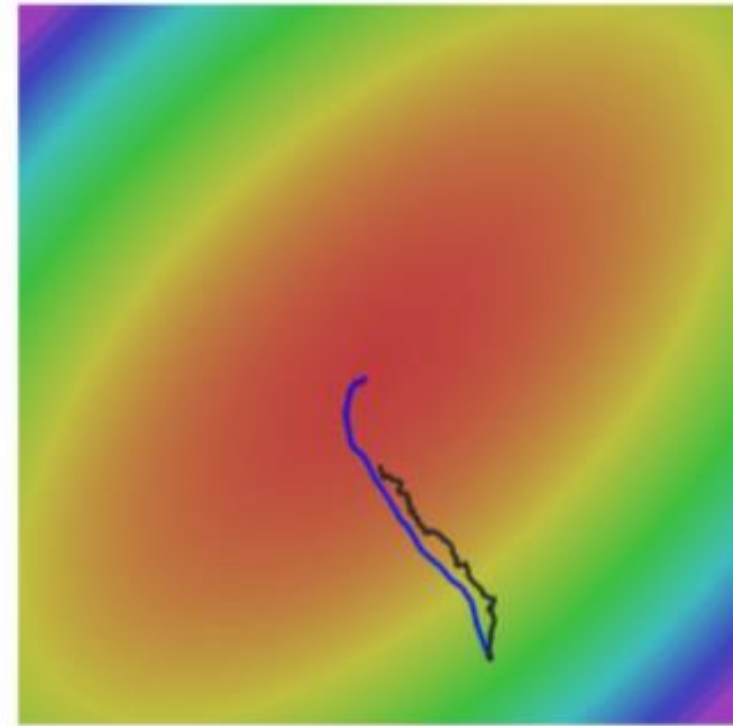
# Optimization (SGD Problem)



**[Saddle Point]**

# Optimization (Momentum)

**[SGD + Momentum]**

Momentum update:

**[Nesterov Momentum]**

Nesterov Momentum



Q. 2가지 방법의 차이점은?

# Optimization (AdaGrad & RMSProp)

### AdaGrad

```python
grad_squared = 0
while True:
    dx = compute_gradient(x)
    grad_squared += dx * dx
    x -= learning_rate * dx / (np.sqrt(grad_squared) + 1e-7)
```

### RMSProp

```python
grad_squared = 0
while True:
    dx = compute_gradient(x)
    grad_squared = decay_rate * grad_squared + (1 - decay_rate) * dx * dx
    x -= learning_rate * dx / (np.sqrt(grad_squared) + 1e-7)
```

# Optimization (Adam)

```
first_moment = 0
second_moment = 0
for t in range(1, num_iterations):
    dx = compute_gradient(x)
    first_moment = beta1 * first_moment  + (1 - beta1) * dx
    second_moment = beta2 * second_moment + (1 - beta2) * dx * dx
    first_unbias = first_moment / (1 - beta1 ** t)
    second_unbias = second_moment / (1 - beta2 ** t)
    x -= learning_rate * first_unbias / (np.sqrt(second_unbias) + 1e-7))
```
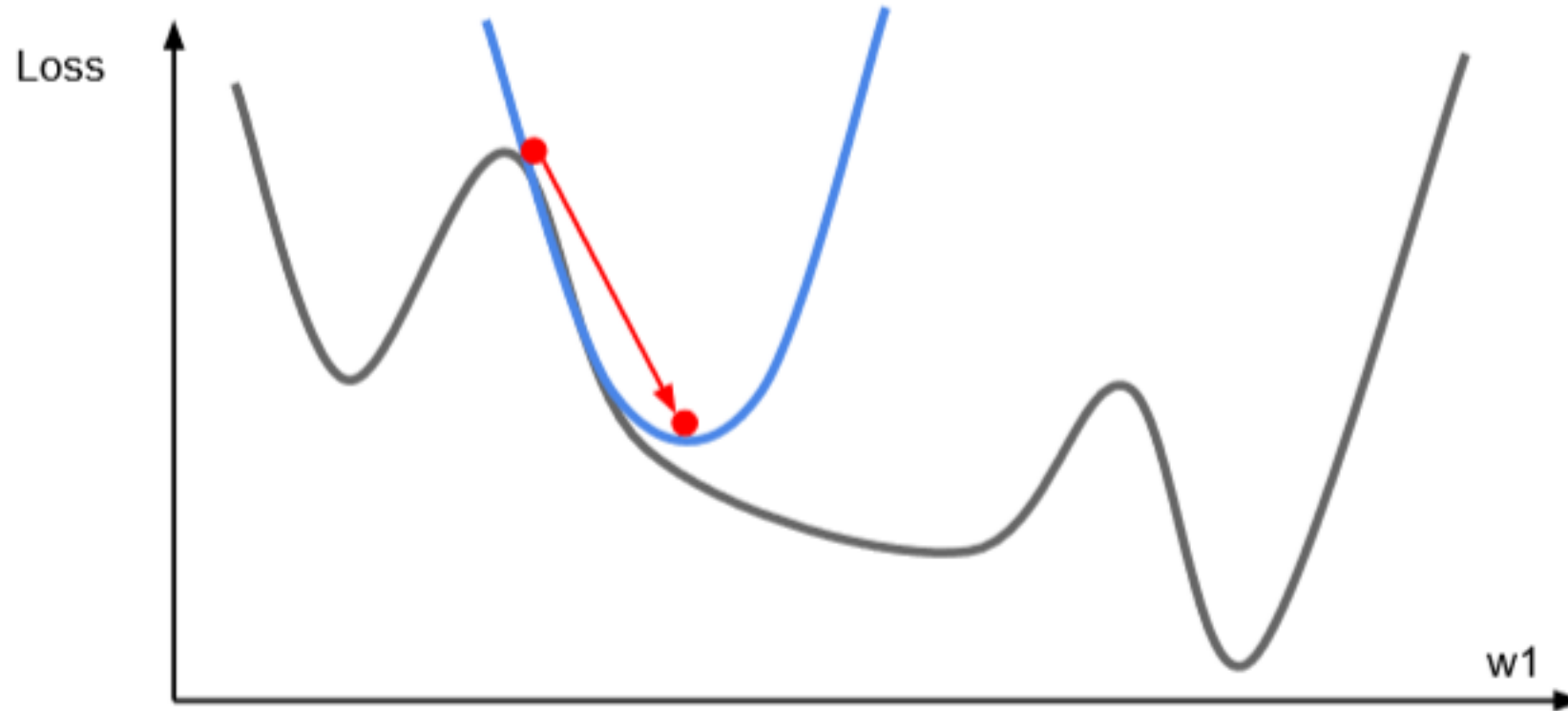
Momentum

Bias correction

AdaGrad / RMSProp

Bias correction for the fact that
first and second moment
estimates start at zero

Adam with beta1 = 0.9,
beta2 = 0.999, and learning_rate = 1e-3 or 5e-4
is a great starting point for many models!

BOAZ

(1)  Use gradient **and Hessian** to form **quadratic** approximation
(2)  Step to the **minima** of the approximation

second-order Taylor expansion:

$$J(\boldsymbol{\theta}) \approx J(\boldsymbol{\theta}_0) + (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^\top \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}_0) + \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}_0)^\top \boldsymbol{H}(\boldsymbol{\theta} - \boldsymbol{\theta}_0)$$

Solving for the critical point we obtain the Newton parameter update:

$$\boldsymbol{\theta}^* = \boldsymbol{\theta}_0 - \boldsymbol{H}^{-1}\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}_0)$$

No hyperparameters!
No learning rate!

- Quasi-Newton methods (**BGFS** most popular): *instead of inverting the Hessian (O(n^3)), approximate inverse Hessian with rank 1 updates over time (O(n^2) each).*

- **L-BFGS** (Limited memory BFGS): *Does not form/store the full inverse Hessian.*
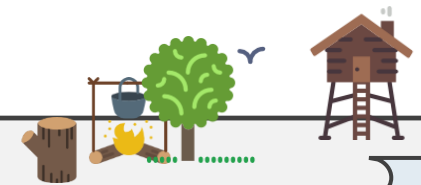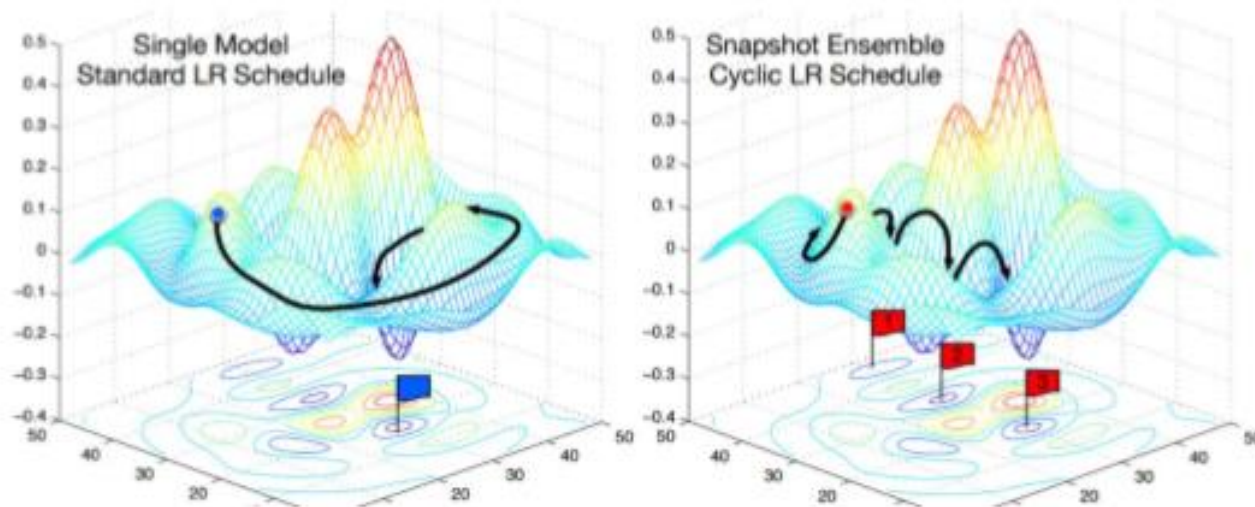
L-BFGS는 Style Transfer와 같이
Full Batch의 size가 작은 경우에 사용됨.

https://pytorch.org/docs/stable/optim.html

1. Train multiple independent models
2. At test time average their results

Enjoy 2% extra performance

# Model Ensembles : Tips and Tricks



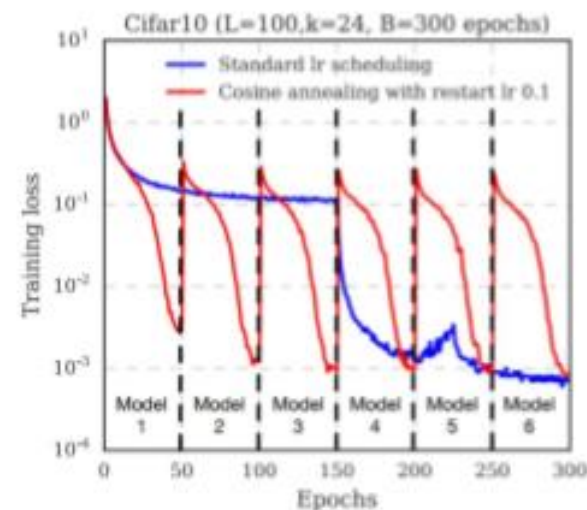Instead of training independent models, use multiple snapshots of a single model during training!

Loshchilov and Hutter, "SGDR: Stochastic gradient descent with restarts", arXiv 2016
Huang et al, "Snapshot ensembles: train 1, get M for free", ICLR 2017
Figures copyright Yixuan Li and Geoff Pleiss, 2017. Reproduced with permission.

Cyclic learning rate schedules can make this work even better!

# Regularization : Add term to loss

$$L = \frac{1}{N} \sum_{i=1}^{N} \sum_{j \neq y_i} \max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + 1) + \boxed{\lambda R(W)}$$

In common use:

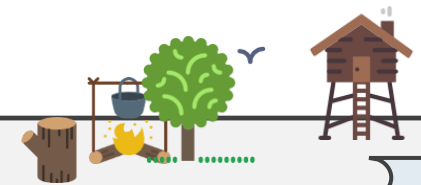**L2 regularization**       $R(W) = \sum_k \sum_l W_{k,l}^2$   (Weight decay)

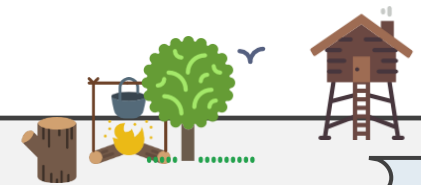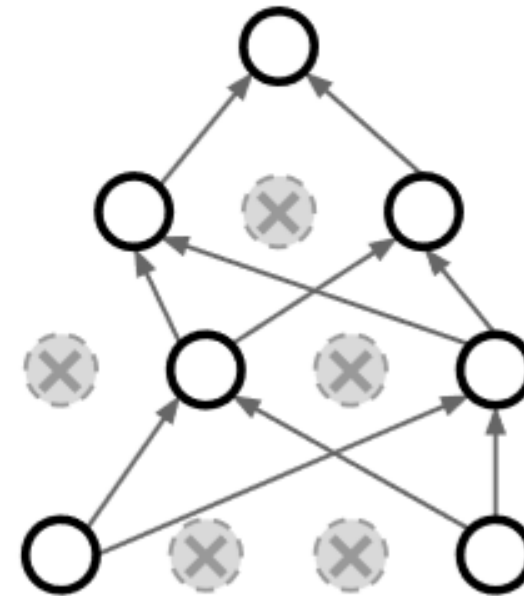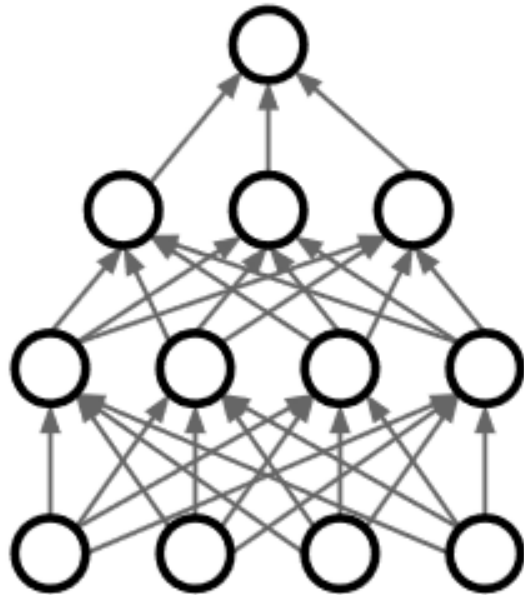L1 regularization       $R(W) = \sum_k \sum_l |W_{k,l}|$

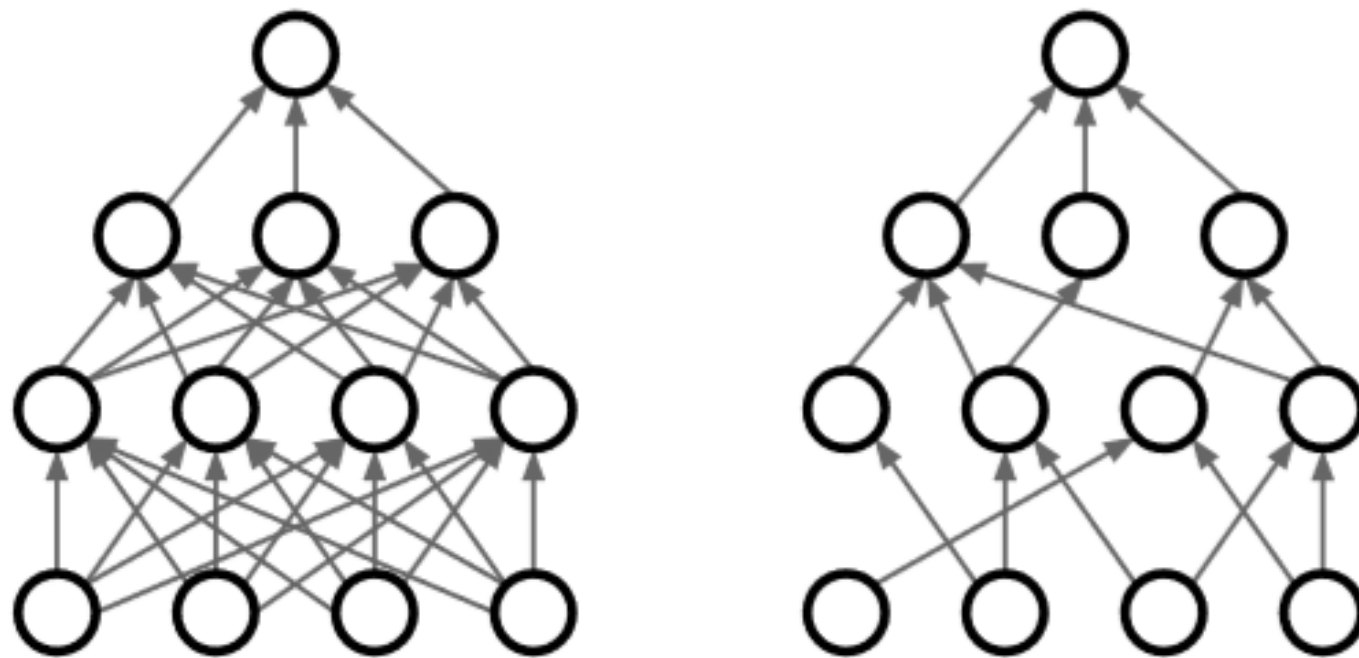Elastic net (L1 + L2)   $R(W) = \sum_k \sum_l \beta W_{k,l}^2 + |W_{k,l}|$

# Regularization : Dropout

In each forward pass, randomly set some neurons to zero
Probability of dropping is a hyperparameter; 0.5 is common

# Regularization : Dropconnect

Q. Dropout과 Dropconnect의 차이점은?

# Regularization : Dropout

```python
p = 0.5 # probability of keeping a unit active. higher = less dropout

def train_step(X):
  """ X contains the data """

  # forward pass for example 3-layer neural network
  H1 = np.maximum(0, np.dot(W1, X) + b1)
  U1 = np.random.rand(*H1.shape) < p # first dropout mask
  H1 *= U1 # drop!
  H2 = np.maximum(0, np.dot(W2, H1) + b2)
  U2 = np.random.rand(*H2.shape) < p # second dropout mask
  H2 *= U2 # drop!
  out = np.dot(W3, H2) + b3

  # backward pass: compute gradients... (not shown)
  # perform parameter update... (not shown)
```

Example forward pass with a 3-layer network using dropout

# Regularization : Dropout

Q. Predict를 할 때, p를 곱해주는 이유는?

```
def predict(X):
  # ensembled forward pass
  H1 = np.maximum(0, np.dot(W1, X) + b1) * p # NOTE: scale the activations
  H2 = np.maximum(0, np.dot(W2, H1) + b2) * p # NOTE: scale the activations
  out = np.dot(W3, H2) + b3
```

At test time all neurons are active always
=> We must scale the activations so that for each neuron:
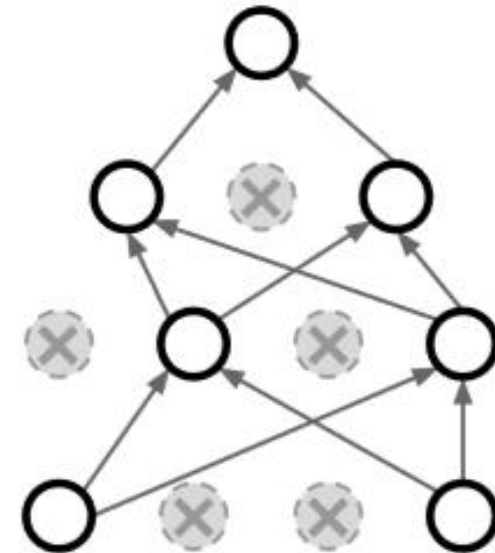output at test time = expected output at training time
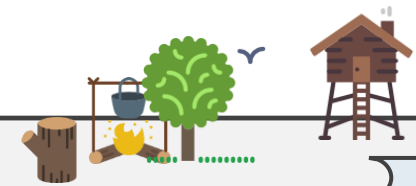
# Regularization : Inverted dropout

```python
p = 0.5 # probability of keeping a unit active. higher = less dropout

def train_step(X):
  # forward pass for example 3-layer neural network
  H1 = np.maximum(0, np.dot(W1, X) + b1)
  U1 = (np.random.rand(*H1.shape) < p) / p # first dropout mask. Notice /p!
  H1 *= U1 # drop!
  H2 = np.maximum(0, np.dot(W2, H1) + b2)
  U2 = (np.random.rand(*H2.shape) < p) / p # second dropout mask. Notice /p!
  H2 *= U2 # drop!
  out = np.dot(W3, H2) + b3

  # backward pass: compute gradients... (not shown)
  # perform parameter update... (not shown)

def predict(X):
  # ensembled forward pass
  H1 = np.maximum(0, np.dot(W1, X) + b1) # no scaling necessary
  H2 = np.maximum(0, np.dot(W2, H1) + b2)
  out = np.dot(W3, H2) + b3
```
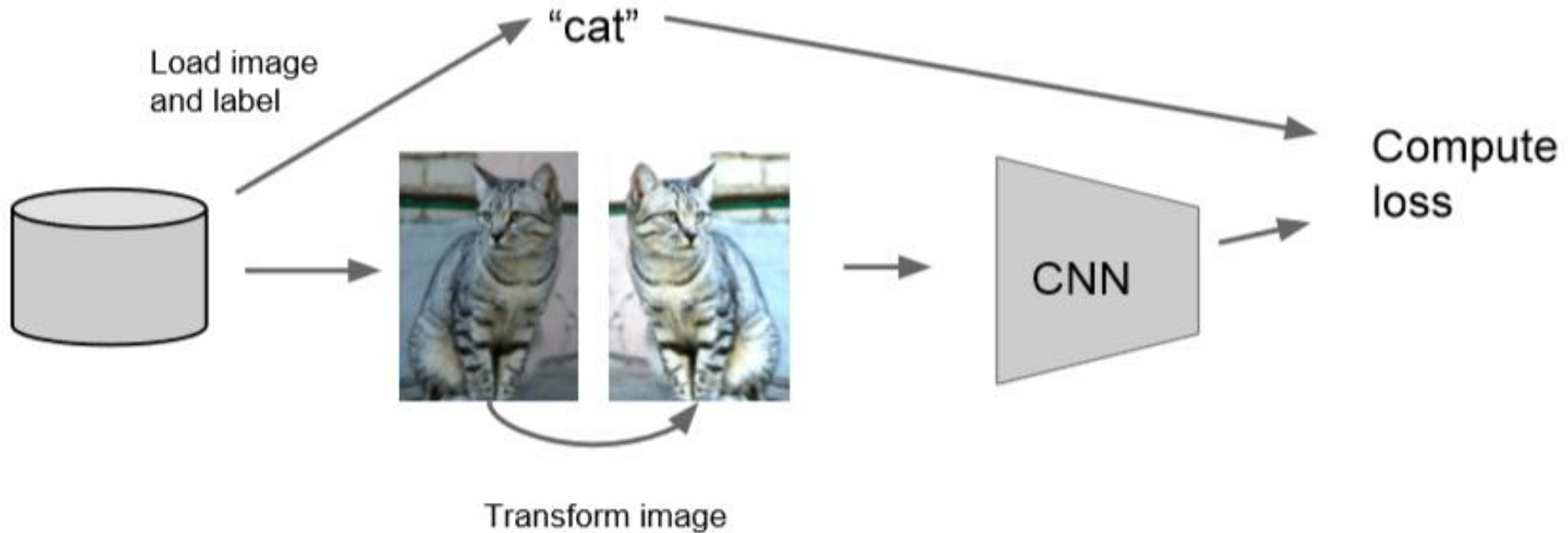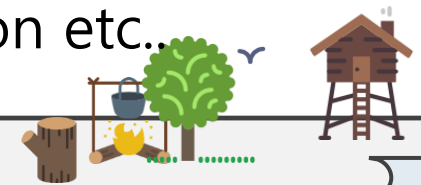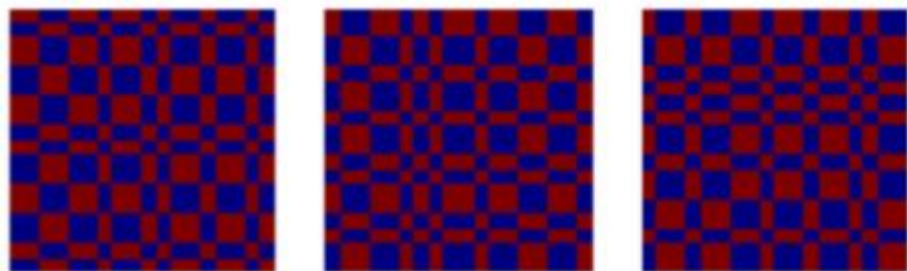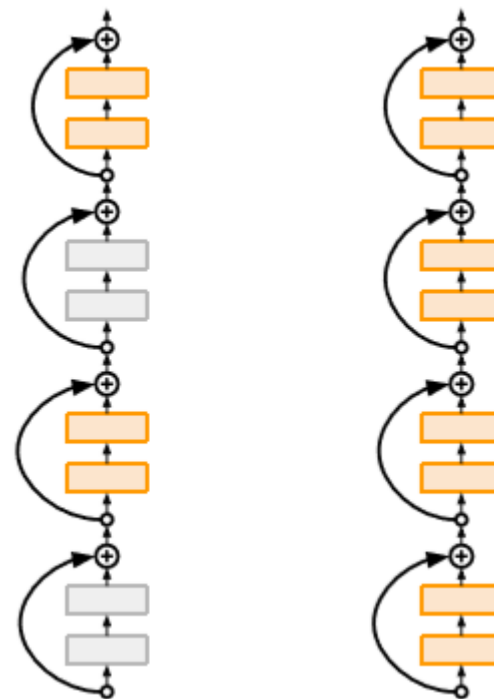
test time is unchanged!

**04**



- Horizontal Flips
- Random crops and scales
- Color Jitter
- Random mix / combinations of rotation, translation etc..

# Regularization : 이외의 방법



[Fractional Max Pooling]

[Stochastic Depth]

# Transfer Learning



## 1. Train on Imagenet

FC-1000
FC-4096
FC-4096

MaxPool
Conv-512
Conv-512

MaxPool
Conv-512
Conv-512

MaxPool
Conv-256
Conv-256

MaxPool
Conv-128
Conv-128

MaxPool
Conv-64
Conv-64

Image

## 2. Small Dataset (C classes)

FC-C
FC-4096
FC-4096

MaxPool
Conv-512
Conv-512

MaxPool
Conv-512
Conv-512

MaxPool
Conv-256
Conv-256

MaxPool
Conv-128
Conv-128

MaxPool
Conv-64
Conv-64

Image

Reinitialize this and train

Freeze these

## 3. Bigger dataset

FC-C
FC-4096
FC-4096

MaxPool
Conv-512
Conv-512

MaxPool
Conv-512
Conv-512

MaxPool
Conv-256
Conv-256

MaxPool
Conv-128
Conv-128

MaxPool
Conv-64
Conv-64

Image

Train these

With bigger dataset, train more layers

Freeze these

Lower learning rate when finetuning; 1/10 of original LR is good starting point

보통 VGG16, ResNet을 이용하여, Transfer Learning

# Transfer Learning



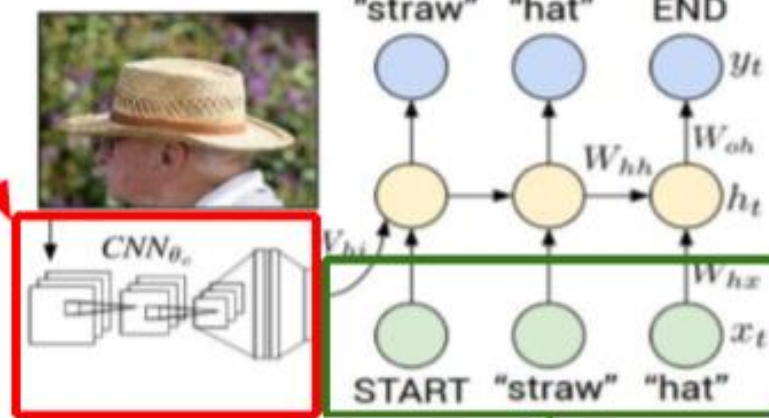|  | very similar dataset | very different dataset |
|---|---|---|
| **very little data** | Use Linear Classifier on top layer | You're in trouble… Try linear classifier from different stages |
| **quite a lot of data** | Finetune a few layers | Finetune a larger number of layers |

# Transfer Learning

# Transfer Learning (Pytorch Model zoo)

| Network | Top-1 error | Top-5 error |
| --- | --- | --- |
| AlexNet | 43.45 | 20.91 |
| VGG-11 | 30.98 | 11.37 |
| VGG-13 | 30.07 | 10.75 |
| VGG-16 | 28.41 | 9.62 |
| VGG-19 | 27.62 | 9.12 |
| VGG-11 with batch normalization | 29.62 | 10.19 |
| VGG-13 with batch normalization | 28.45 | 9.63 |
| VGG-16 with batch normalization | 26.63 | 8.50 |
| VGG-19 with batch normalization | 25.76 | 8.15 |
| ResNet-18 | 30.24 | 10.92 |
| ResNet-34 | 26.70 | 8.58 |
| ResNet-50 | 23.85 | 7.13 |
| ResNet-101 | 22.63 | 6.44 |

https://pytorch.org/docs/stable/torchvision/models.html

**End** **참고 자료**

CS231n : http://cs231n.stanford.edu/syllabus.html

Pytorch Optimization : https://pytorch.org/docs/stable/optim.html

PyTorch Model Zoo : https://pytorch.org/docs/stable/torchvision/models.html