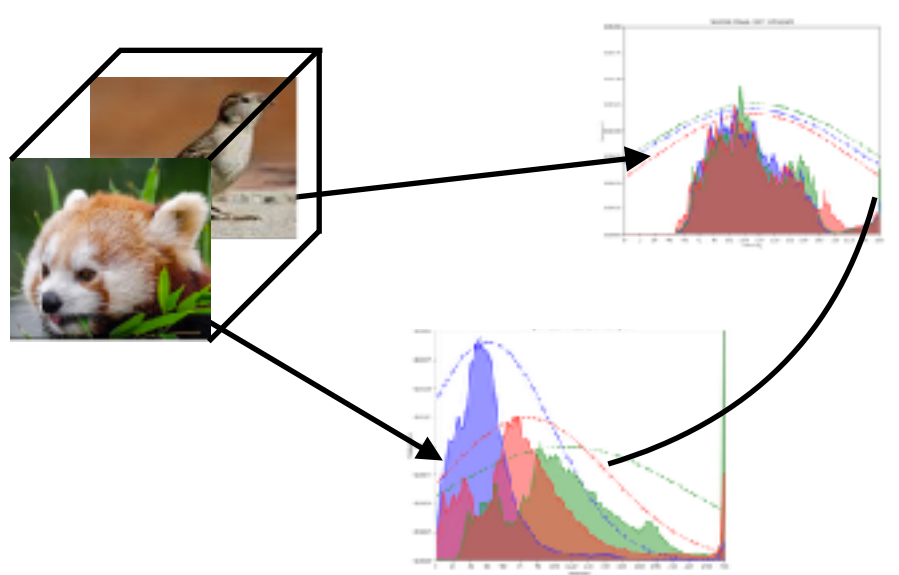


Batch Normalization

Why normalization need?



Problem

Different pixel distribution make hard train image features.
To avoid , gradient vanishing , gradient exploding , set learning rate to low

Solution

Normalization
Normalization

$$\frac{X - X_{\min}}{X_{\max} - X_{\min}}$$

Batch normalization

Batch normalization

STEP 1

$$x'_{\text{each_pixel}} = \frac{x_{\text{each_pixel}} - E(x)_{\text{batch}}}{\text{var}(x)_{\text{batch}}}$$

Neural network 을 학습

STEP 2

$$X = \gamma * X + \beta$$

"Note that simply normalizing each input of a layer may change what the layer can represent. For instance, normalizing the inputs of a sigmoid would constrain them to the linear regime of the nonlinearity."

STEP 3

$$\frac{\sum_i E(x)_{i,\text{batch}}}{K} = E(x)_{\text{test}}$$

$$\frac{\sum_i \text{var}(x)_{i,\text{batch}}}{K} = \text{var}(x)_{\text{test}}$$

Test 학습을 할때 마다 mean. , var 이 달라지기 때문에 변하지 않는 mean , var 이 필요하다.

성능 평가

```
import tensorflow as tf
x = tf.constant(10.0, name='input')
w = tf.Variable(0.8, name='weight')
y = tf.multiply(w, x, name='output')
y_ = tf.constant(0.0, name='correct_value')
loss = tf.pow(y - y_, 2, name='loss')
train_step = tf.train.GradientDescentOptimizer(0.025).minimize(loss)

for value in [x, w, y, y_, loss]:
    tf.summary.scalar(value.op.name, value)

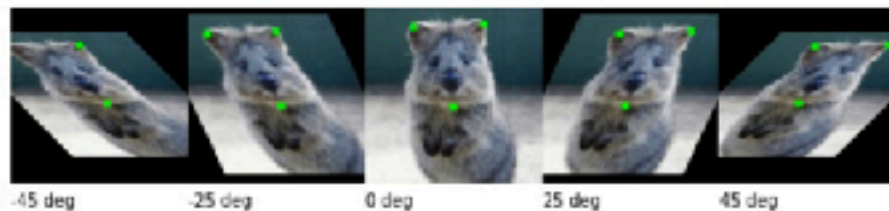
summaries = tf.summary.merge_all()

sess = tf.Session()
summary_writer = tf.summary.FileWriter('log_simple_stats', sess.graph)

sess.run(tf.initialize_all_variables())
for i in range(100):
    if i % 10 == 0:
        print("epoch {}, output: {}".format(i, sess.run(y)))
        summary_writer.add_summary(sess.run(summaries), i)
    sess.run(train_step)
```


AUGMENTATION

Affine: Shear



Affine: Modes



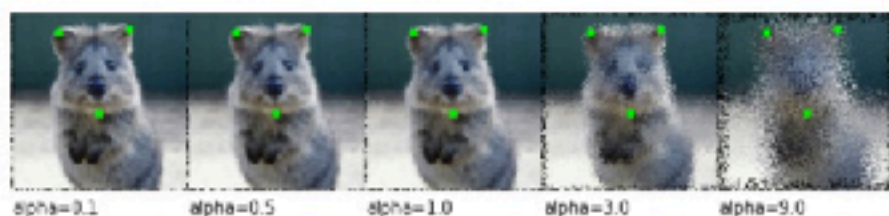
Affine: cval



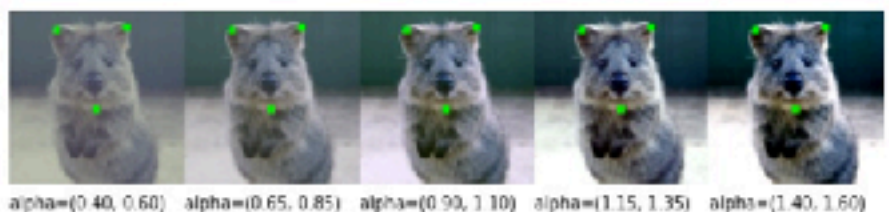
Affine: all



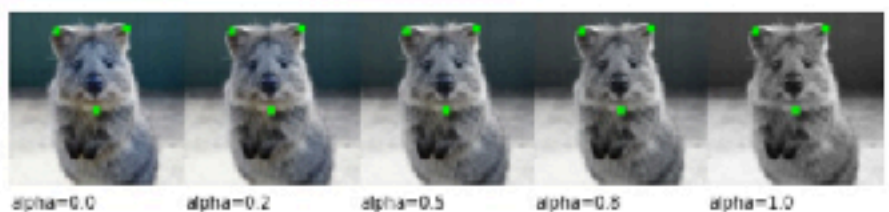
ElasticTransformation
(sigma=0.2)



ContrastNormalization
(per channel)



Grayscale



GaussianBlur



AdditiveGaussianNoise



AdditiveGaussianNoise
(per channel)



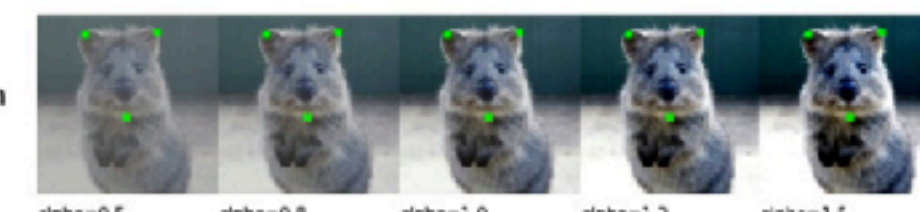
Dropout



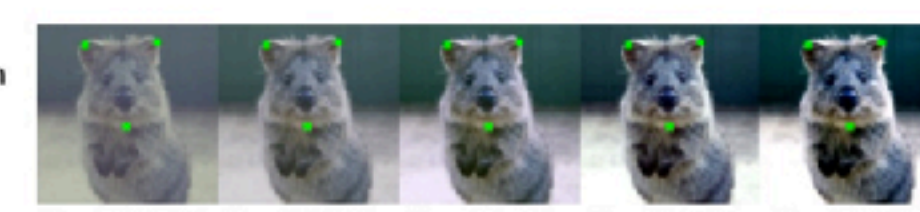
Dropout
(per channel)



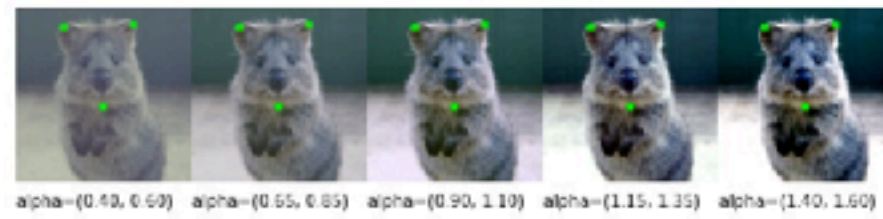
ContrastNormalization



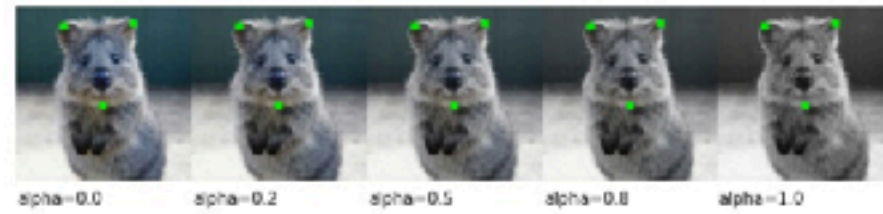
ContrastNormalization
(per channel)



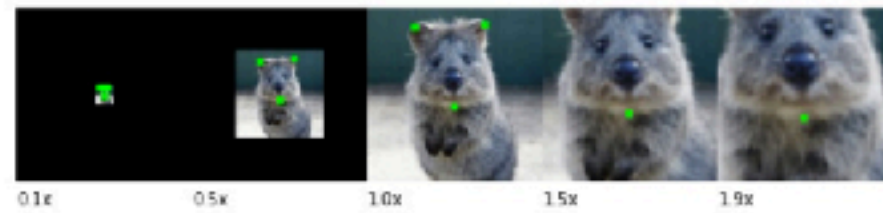
Contrast Normalization
(per channel)



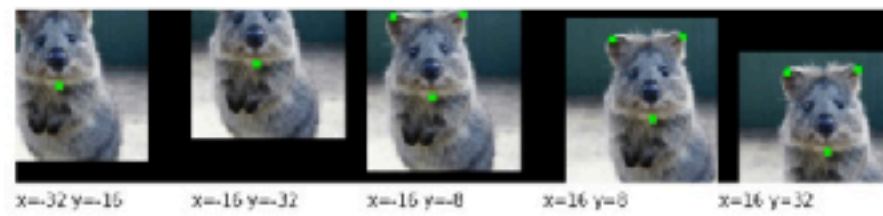
Grayscale



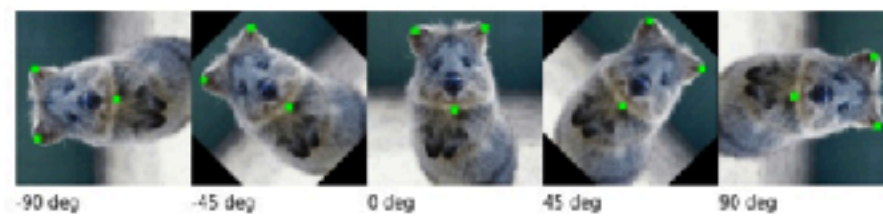
Affine: Scale



Affine: Translate



Affine: Rotate



Gradient Optimization

<http://sebastianruder.com/optimizing-gradient-descent/>

모든 파라미터 θ :parameter 을 최소화 하기 위해

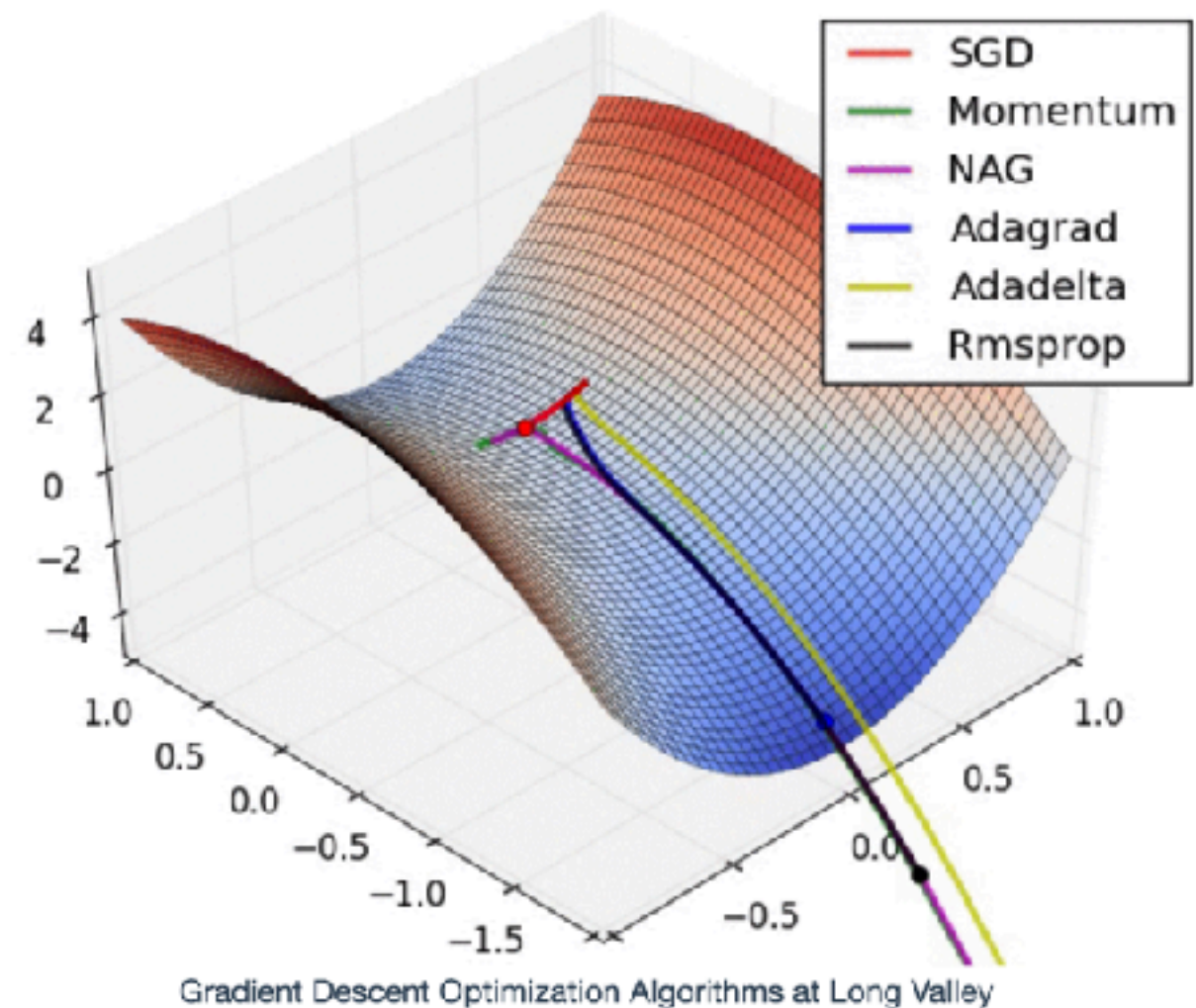
$J(\theta)$:Loss function 을 줄이기 위해 .

$\nabla_{\theta}J(\theta)$ 을 이용한다.

↓
어떻게?

↓
기울기의 반대방향으로

↓
$$\theta = \theta - \eta \nabla_{\theta}J(\theta)$$



Optimizer

Momentum

SGD(Stochastic Gradient Descent (SGD))

NAG(Nesterov Accelerated Gradient)

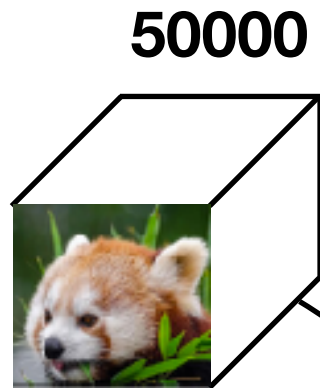
Adagrad

RMSProp

AdaDelta

Adam

SGD



BATCH DESCENT GRADIENT

한장을 꺼내서 한번 학습하고,
차례대로 50000장의 이미지를 모두 꺼내 학습한다

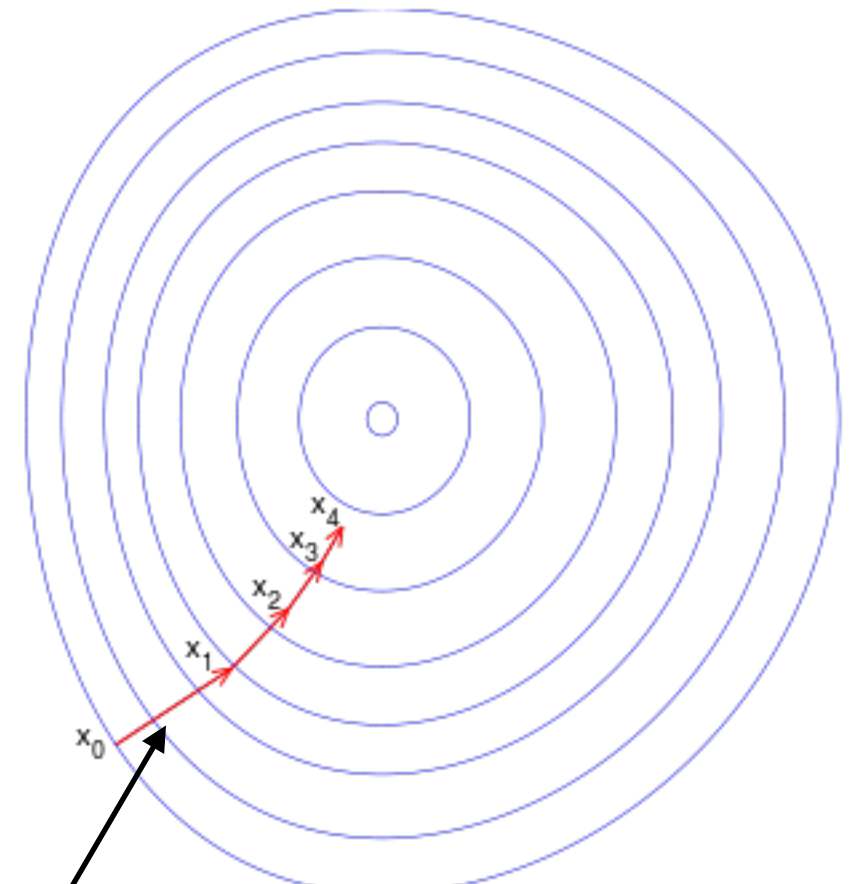
$$\theta = \theta - \eta \nabla_{\theta} J(\theta)$$



η : learning rate

$$\theta_{batch} = \theta_{batch} - \eta \nabla_{\theta} J(\theta)_{batch}$$
$$[J(\theta)_{batch} = \sum J(\theta)_{batch} / batch]$$

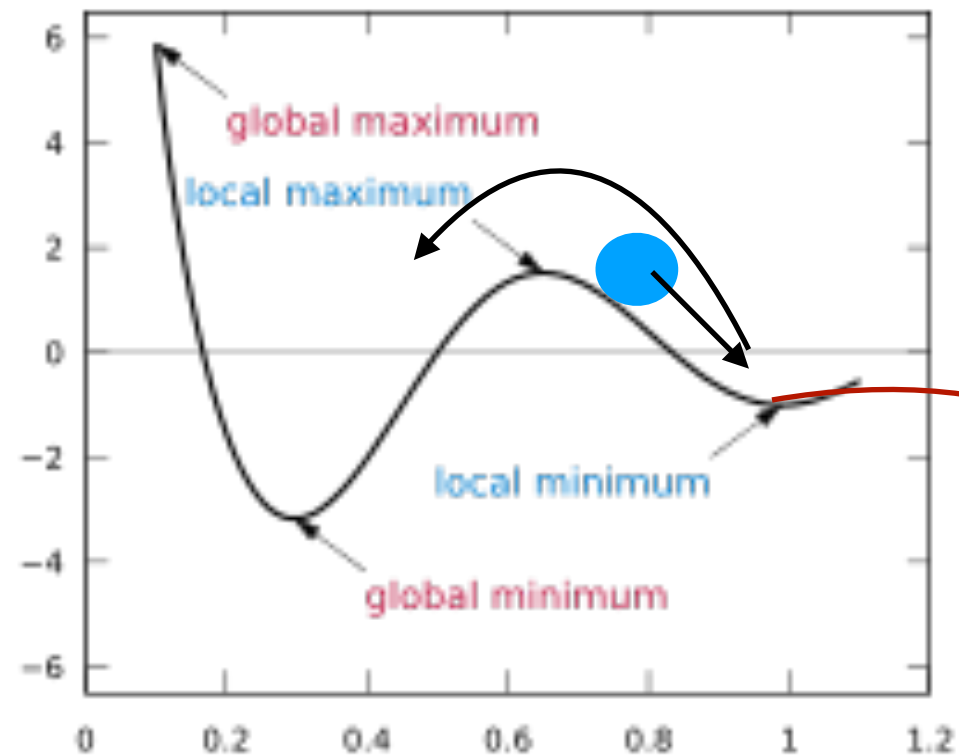
단점.
이동속도가 현저하게 느리다
Local minima에 쉽게 빠진다



등고선이 크기(기울기) 일정하다 라고 생각하기 때문에
SGD는 optimization 하는데 속도가 느리고 한계점이 있다.

그러면 기울기에 따라 변화량(∇)에 따라 learning rate 을 변하게 하면 안될까?

Local minima



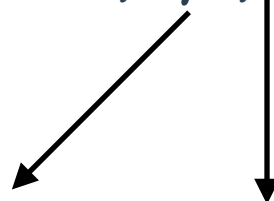
$$\theta_{batch} = \theta_{batch} - \eta \nabla_{\theta} J(\theta)_{batch}$$

$$\nabla_{\theta} J(\theta)_{batch} = 0$$

학습 종료

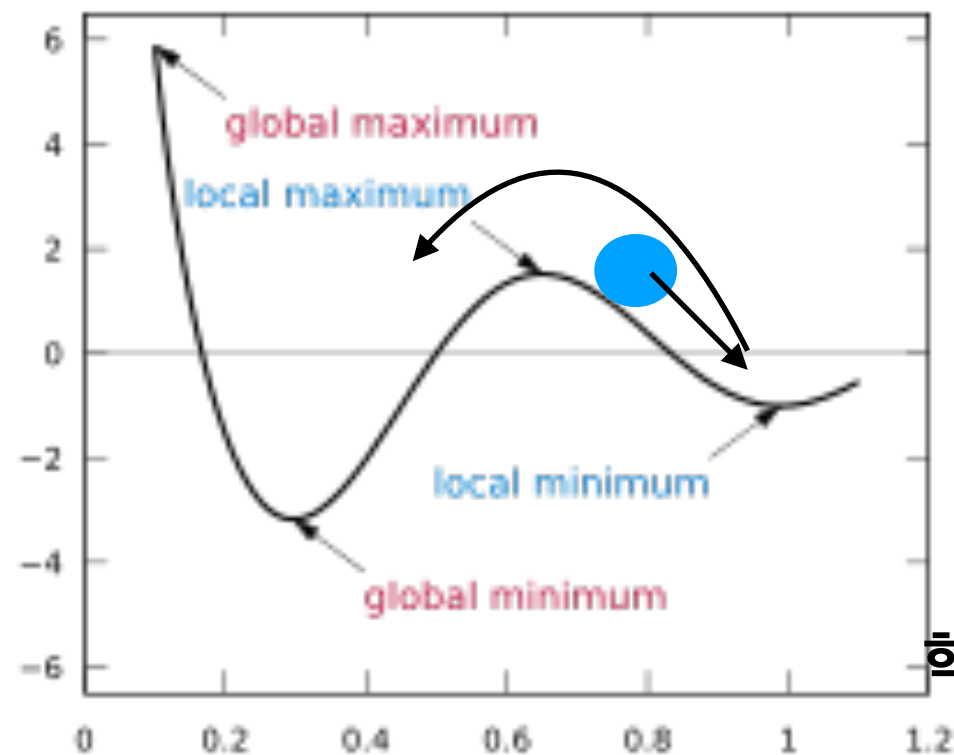
Momentum

Learning rate 을 고정하지 말고 변화량에 따라 learning rate 을 변화시키자!

$$v_t = \gamma v_{t-1} + \eta \nabla_{\theta} J(\theta)$$


γ : momentum

$$v_t = \eta \nabla_{\theta} J(\theta)_t + \gamma \eta \nabla_{\theta} J(\theta)_{t-1} + \gamma^2 \eta \nabla_{\theta} J(\theta)_{t-2} + \dots$$



이 공이 내려올때 더 큰 가중치를 받고 올라가면
Local minima을 뛰어 넘을수 있지 않을까?

학습 속도가 느리다.

<https://github.com/SoulDuck/resnet>

Momentum

학습 속도가 느리다.

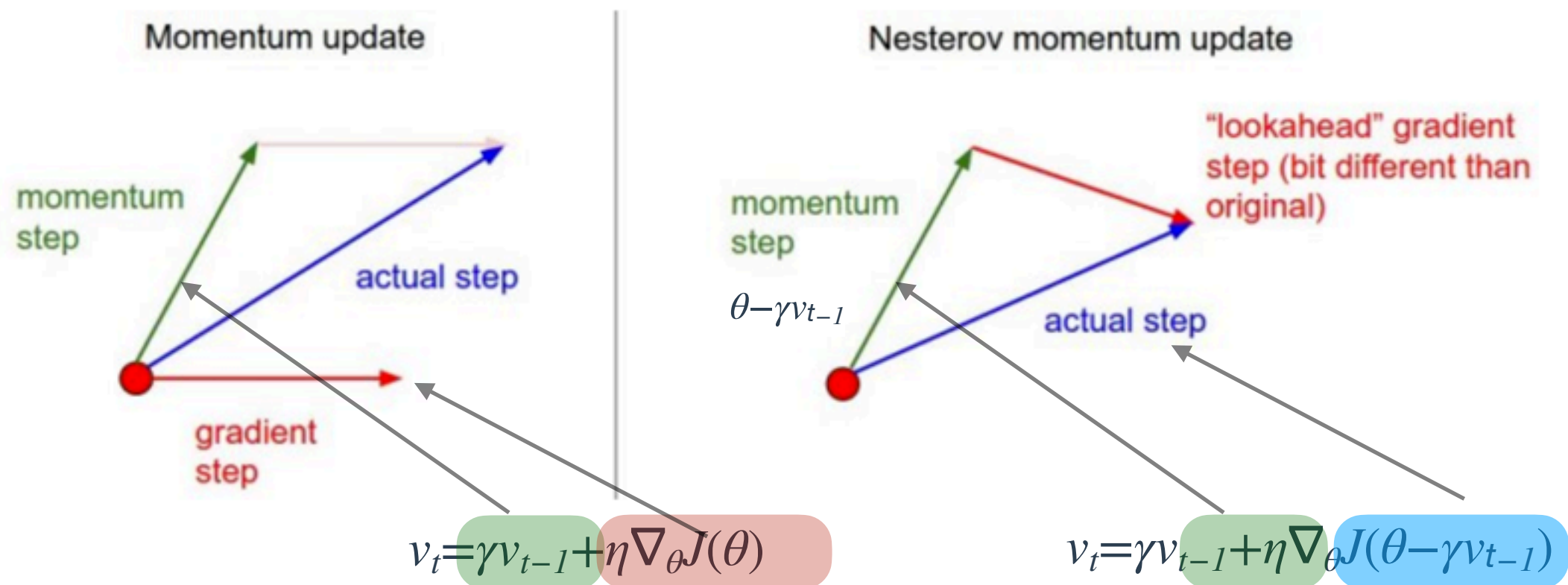
변수를 저장해야 한다. *Because... $v_t = \eta \nabla_{\theta} J(\theta)_t + \gamma \eta \nabla_{\theta} J(\theta)_{t-1} + \gamma^2 \eta \nabla_{\theta} J(\theta)_{t-2} + \dots$*

Global minimum을 빠져나간다.

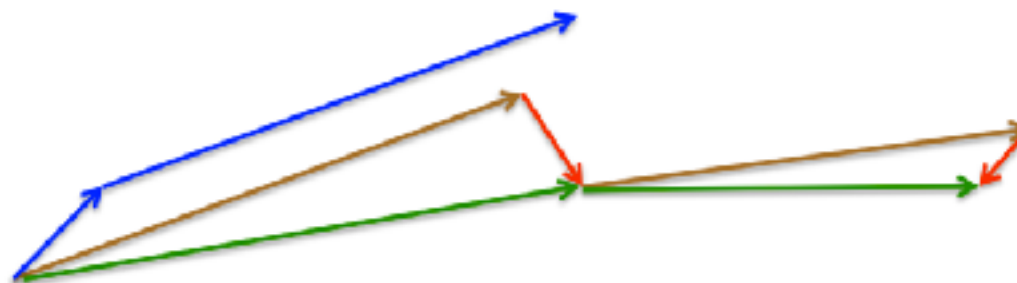


멈출때 관성터를 적게 줘서 멈추게 하는 방법이 있을까?

Nesterov Accelerated Gradient



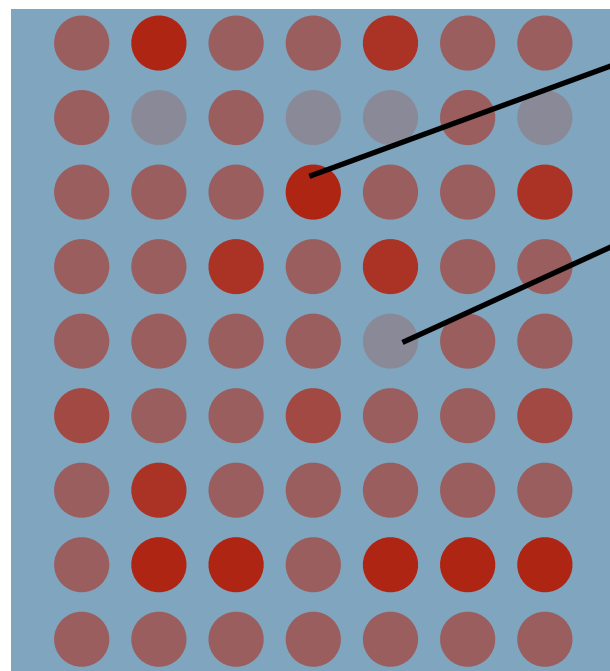
momentum이 공이 흘러내리는 방향과 다른쪽의 힘을 준다...



Adagrad (Adaptive gradient)

Adaptive:

이미 많이 변화한 변수들은 변화량을 적게하고 적게 변화한 변수들은 변화량을 크게하자



변화량이 큰 변수 G_{18}

step size decay

변화량이 작은 변수 G_{33}

step size을 신경쓰지 않는다.

$$G_t = G_{t-1} + (\nabla_{\theta} J(\theta_t))^2$$

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{G_t + \epsilon}} \cdot \nabla_{\theta} J(\theta_t)$$

```
for i in range(10000):
    msg = '\r Progress {0}/{1}'.format(i, 10000)
    sys.stdout.write(msg)
    sys.stdout.flush()
    if i < 1000:
        lrn_point = 0.01
    elif i < 2000:
        lrn_point = 0.001
    elif i < 5000:
        lrn_point = 0.0001
    else:
        lrn_point = 0.0001
```

Adaptive의 가장큰 단점

- 그래디언트의 축적으로 인해 param이 변하지 않는다.



- 그래디언트의 축적을 제한 하는 방법?

Adadelta

The diagram illustrates the derivation of the Adadelta update rule. It starts with the standard update rule $\theta_{t+1} = \theta_t + \Delta\theta_t$, where $\Delta\theta_t$ is highlighted in a blue box. An arrow points from this box to a light blue box containing the formula $\Delta\theta_t = -\frac{RMS[\Delta\theta]_{t-1}}{RMS[g]_t} g_t$, with $RMS[\Delta\theta]_{t-1}$ highlighted in a smaller blue box. From this box, an arrow points to the right, leading to a sequence of equations: $\Delta\theta_t = -\frac{\eta}{\sqrt{E[g^2]_t + \epsilon}} g_t$, followed by $RMS[\Delta\theta]_t = \sqrt{E[\Delta\theta^2]_t + \epsilon}$, then $\Delta\theta_t = -\frac{RMS[\Delta\theta]_{t-1}}{RMS[g]_t} g_t$, and finally $\theta_{t+1} = \theta_t + \Delta\theta_t$. A vertical line connects the first two equations, and another vertical line connects the third and fourth. A downward arrow points from the final equation to the text below.

$$\theta_{t+1} = \theta_t + \Delta\theta_t$$
$$\Delta\theta_t = -\frac{RMS[\Delta\theta]_{t-1}}{RMS[g]_t} g_t$$
$$\Delta\theta_t = -\frac{\eta}{\sqrt{E[g^2]_t + \epsilon}} g_t$$
$$RMS[\Delta\theta]_t = \sqrt{E[\Delta\theta^2]_t + \epsilon}$$
$$\Delta\theta_t = -\frac{RMS[\Delta\theta]_{t-1}}{RMS[g]_t} g_t$$
$$\theta_{t+1} = \theta_t + \Delta\theta_t$$

1.learning rate term이 없어서 정할 필요가 없다.

RMS Prop

$$E[g^2]_t = 0.9E[g^2]_{t-1} + 0.1g_t^2$$
$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{E[g^2]_t + \epsilon}} g_t$$

힐튼이 제시한 Optimizier

계산량이 적고 수식이 복잡하지 않아 많이 사용한다. 성능도 합리적이다.

기본적으로 추천하는 learning rate = 0.0001

ADAM

Adaptive Moment Estimate

RMSProp에다가 momentum 값을 더하자

RMS

$$E[g^2]_t = 0.9E[g^2]_{t-1} + 0.1g_t^2$$
$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{E[g^2]_t + \epsilon}} g_t$$

ADAM

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\hat{v}_t + \epsilon}} \hat{m}_t.$$

Momentum

$$v_t = \gamma v_{t-1} + \eta \nabla_{\theta} J(\theta)$$

Adam

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$$
$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$$

$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}$$

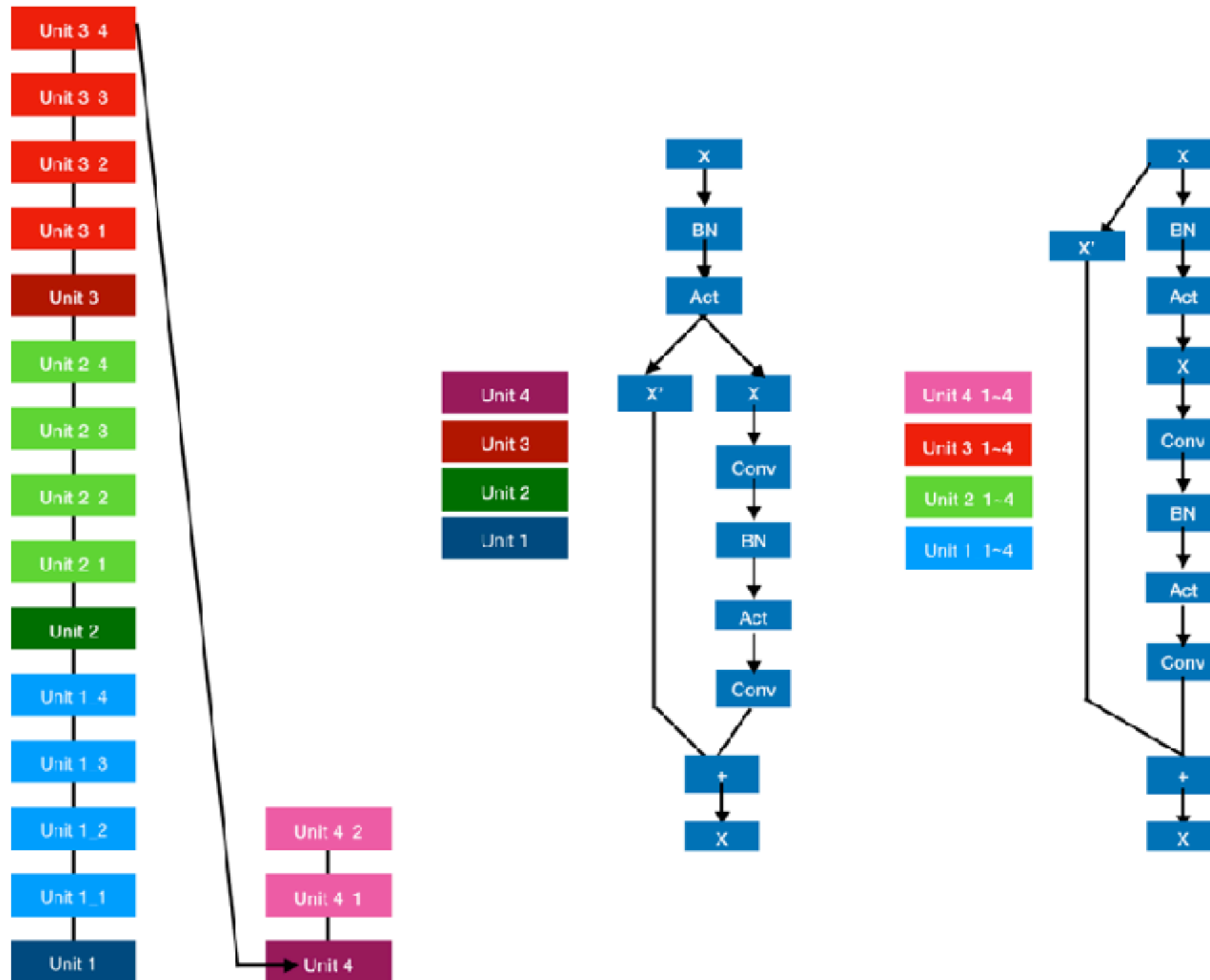
$$\hat{v}_t = \frac{v_t}{1 - \beta_2^t}$$

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{\hat{v}_t} + \epsilon} \hat{m}_t.$$

Like momentum , this sentence remember
t-1's **Gradient**

To avoid over number of optimizer ,
Using m_t^{\wedge}

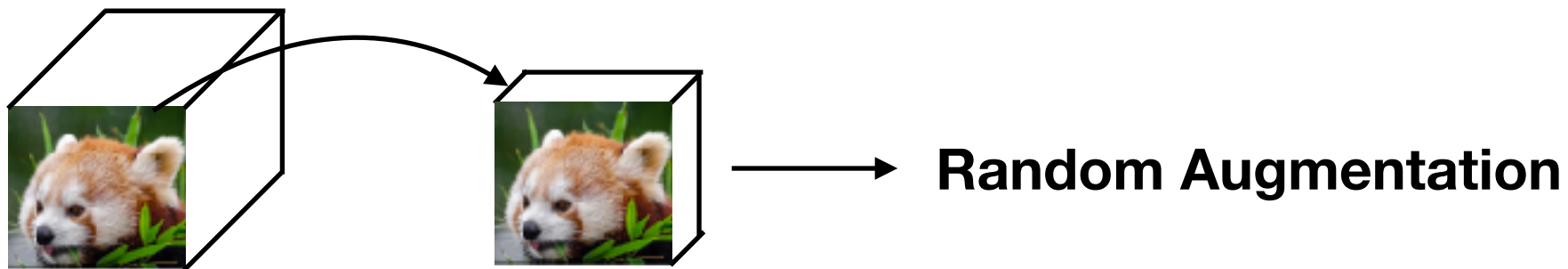
Wide resnet을 이용한 실습



Tensor flow

Batch Images Random Augmentation

```
image = tf.image.resize_image_with_crop_or_pad(image , image_size+4 , image_size+4)
image = tf.random_crop(image , [image_size , image_size ,3])
image = tf.image.random_flip_left_right(image)
image = tf.image.random_flip_up_down(image)
# Brightness / saturatio / constrast provides samll gains 2%~5% on cifar
image = tf.image.random_brightness(image , max_delta=63. / 255.)
image = tf.image.random_saturation(image , lower=0.5 , upper=1.8)
image = tf.image.per_image_standardization(image)
```



Wide resnet을 이용한 실습

- Random batch augmentation
- SGD , Momentum Optimizer
- Tensor board 을 이용한 Visualization