

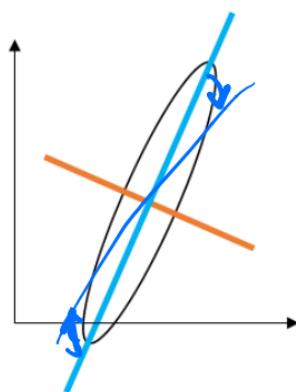


# Lecture 8. Training Neural Networks III

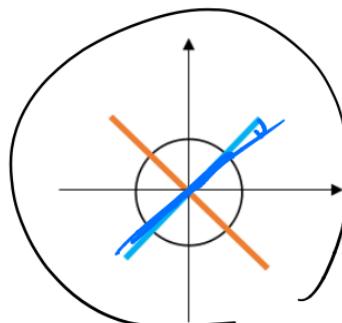
# Review

## Data Normalization.

Before normalization



After normalization



왜?

어떻게?

$$x_{\text{new}} = \frac{x - x_{\min}}{x_{\max} - x_{\min}} \quad \text{Data} \rightarrow [0, 1]$$

$$x_{\text{new}} = \frac{x - \mu}{\sigma}$$

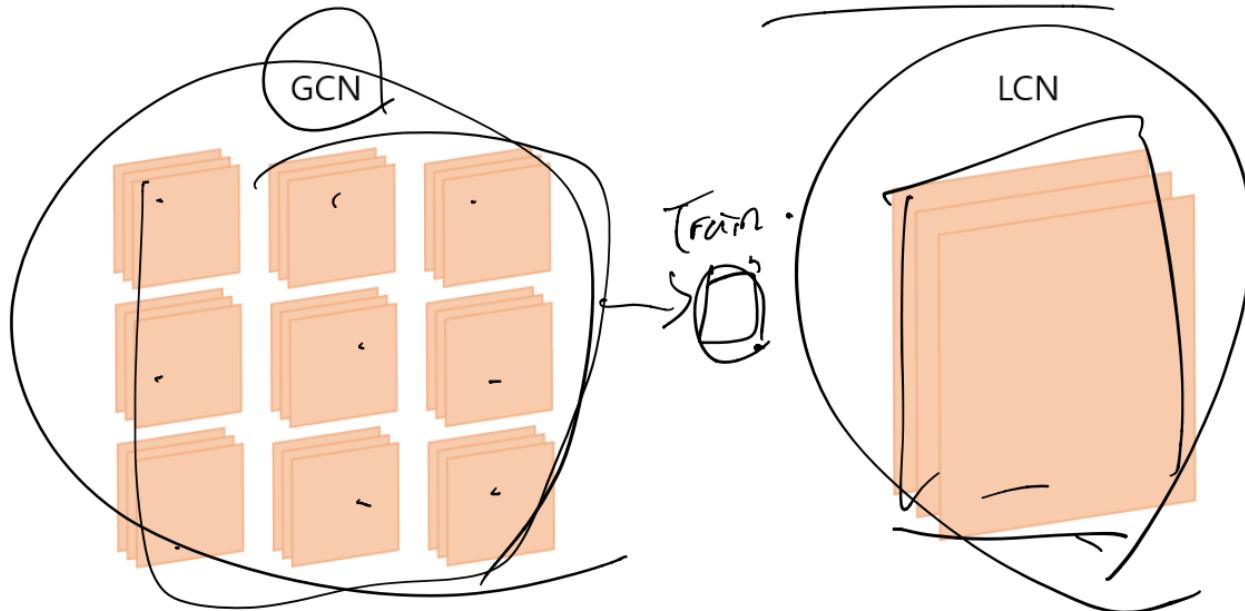
•  $x \sim N(0, 1)$ .

zero-centered

Sigmoid./ReLU

# Review

Image. 모든 이미지의 대화형 GCN  
각각 이미지의 "LCN"



# Review

## Batch Normalization

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_1 \dots m\}$ ;

Parameters to be learned:  $\gamma, \beta$

**Output:**  $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{mini-batch mean}$$

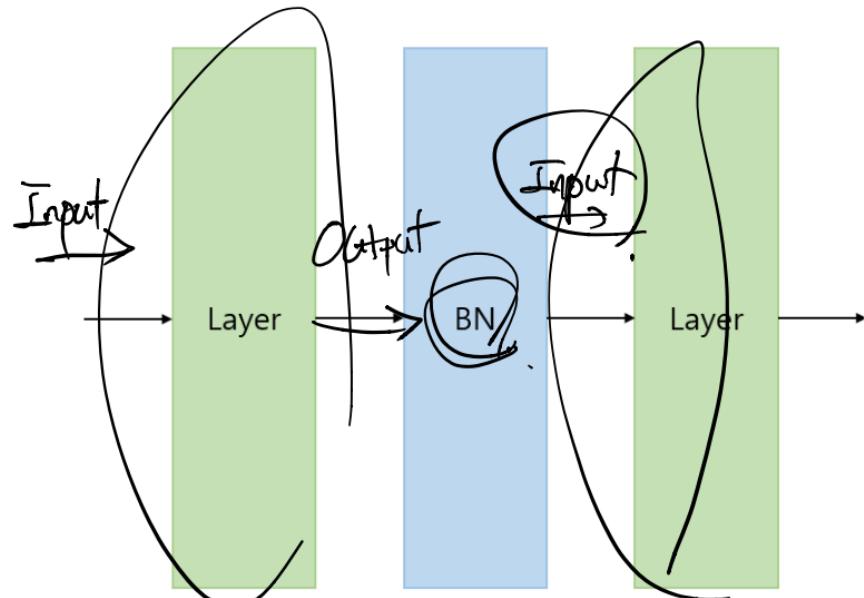
$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{scale and shift}$$

$\gamma, \beta$ .

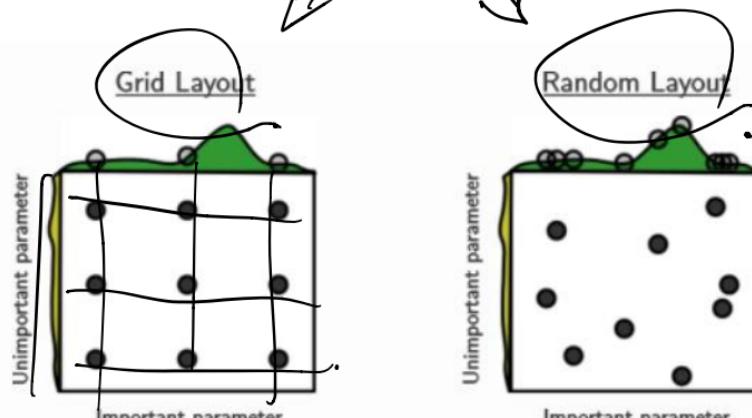
이재현, 김동주



# Review

Hyperparameter search.

+ Auto ML

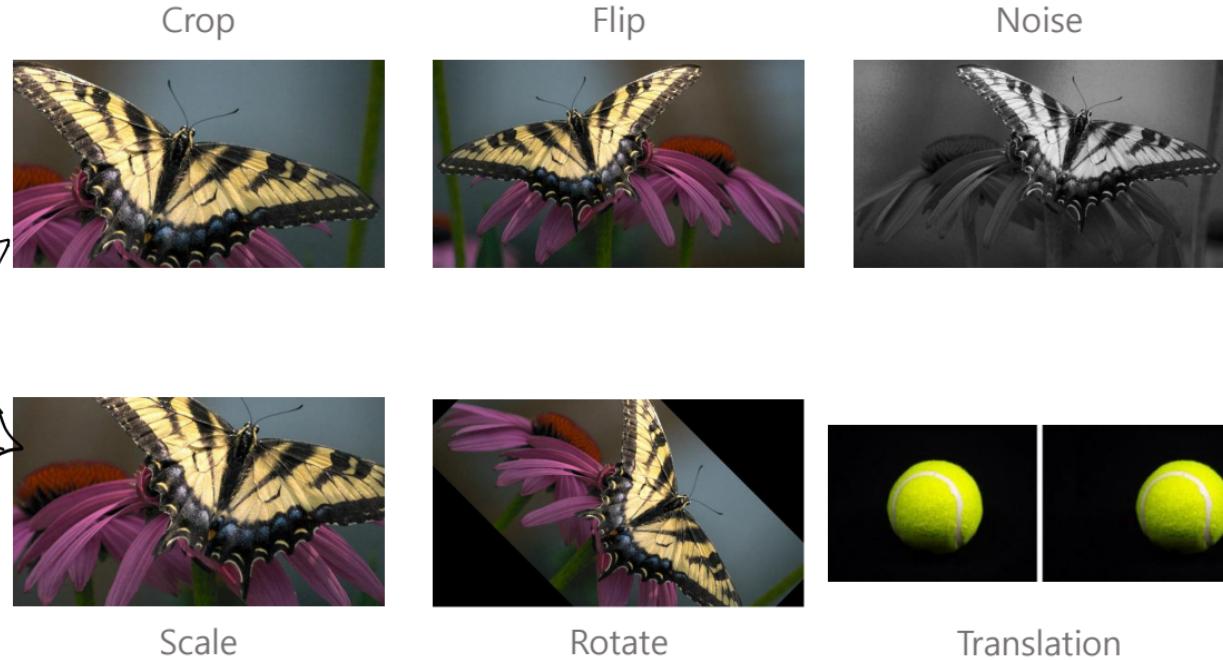
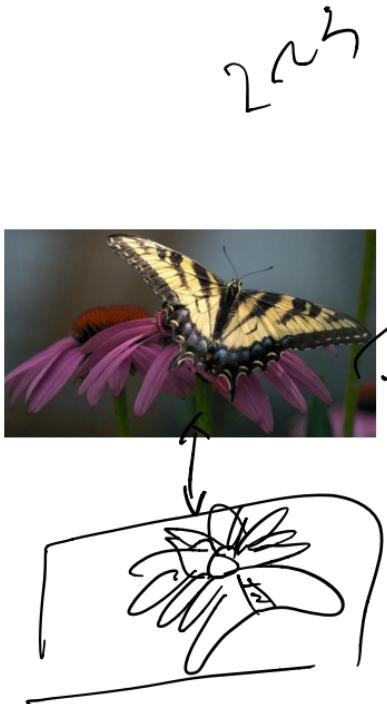


Grid Search 결과와 Random Search 결과 비교 예시

[Bergstra and Bengio(2012)]

# Review

## *Image augmentation*



# Today's Contents

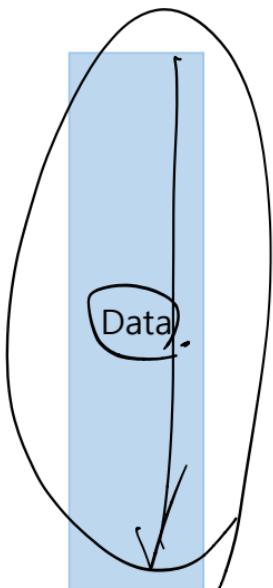
## 1. Change Optimization Process

- Gradient based method - Optimizer
- Regularization - Dropout  
*(f Ensemble)*

# Gradient Descent Optimization Method

1. Gradient Descent
2. Stochastic Gradient Descent (SGD)
3. Momentum
4. Nesterov Accelerated Gradient (NAG)
5. Adagrad
6. RMSProp
7. AdaDelta
8. Adam

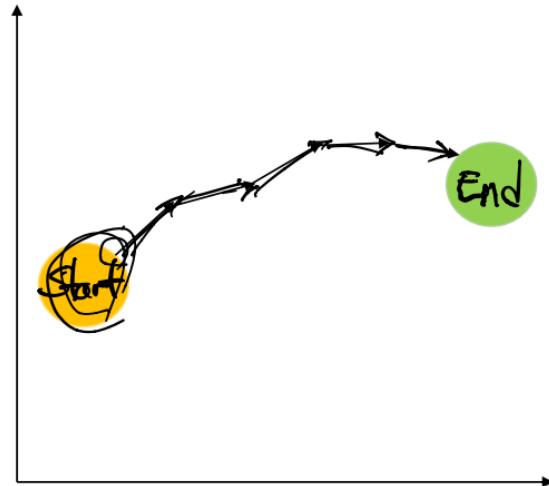
# Gradient Descent (GD)



$$w_{\text{new}} = w_{\text{old}} - \eta \times \frac{\partial L}{\partial w}$$

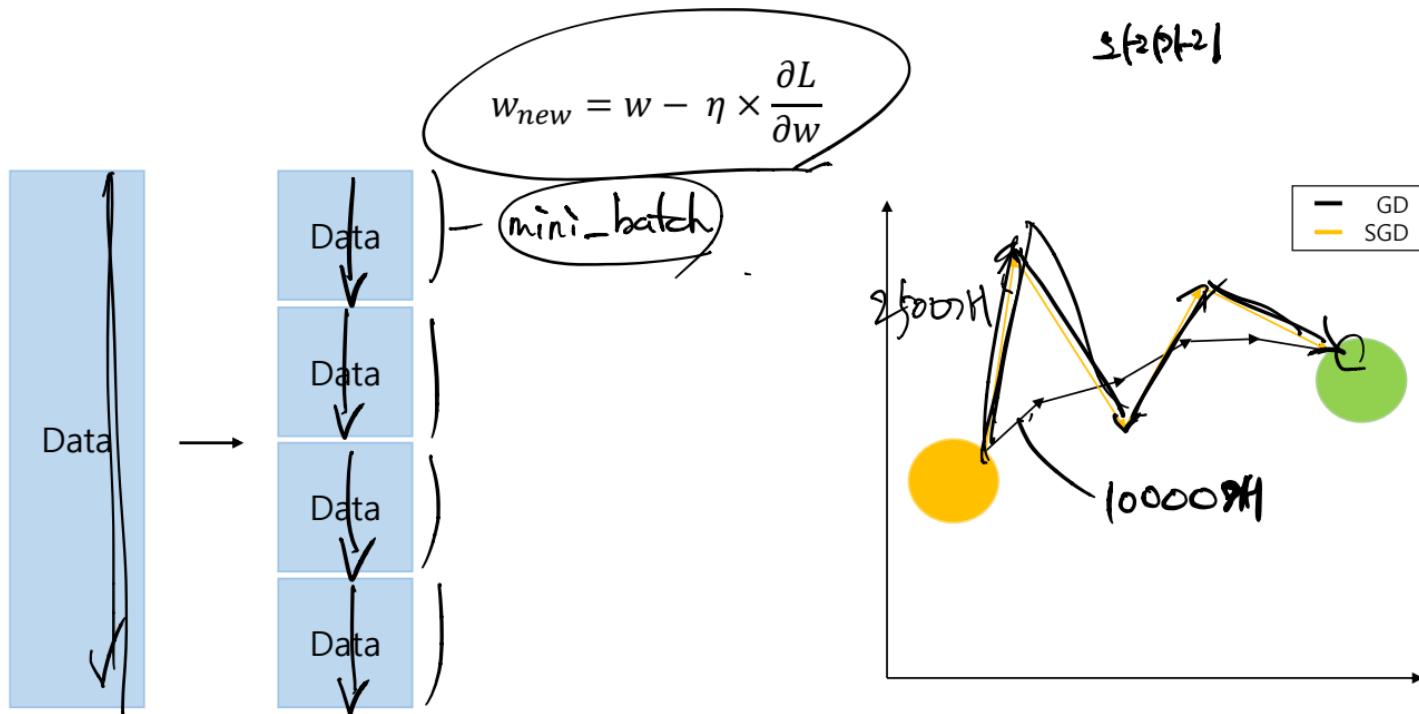
gradient  
learning rate

거의 정착하다!  
시간이 너무 오래

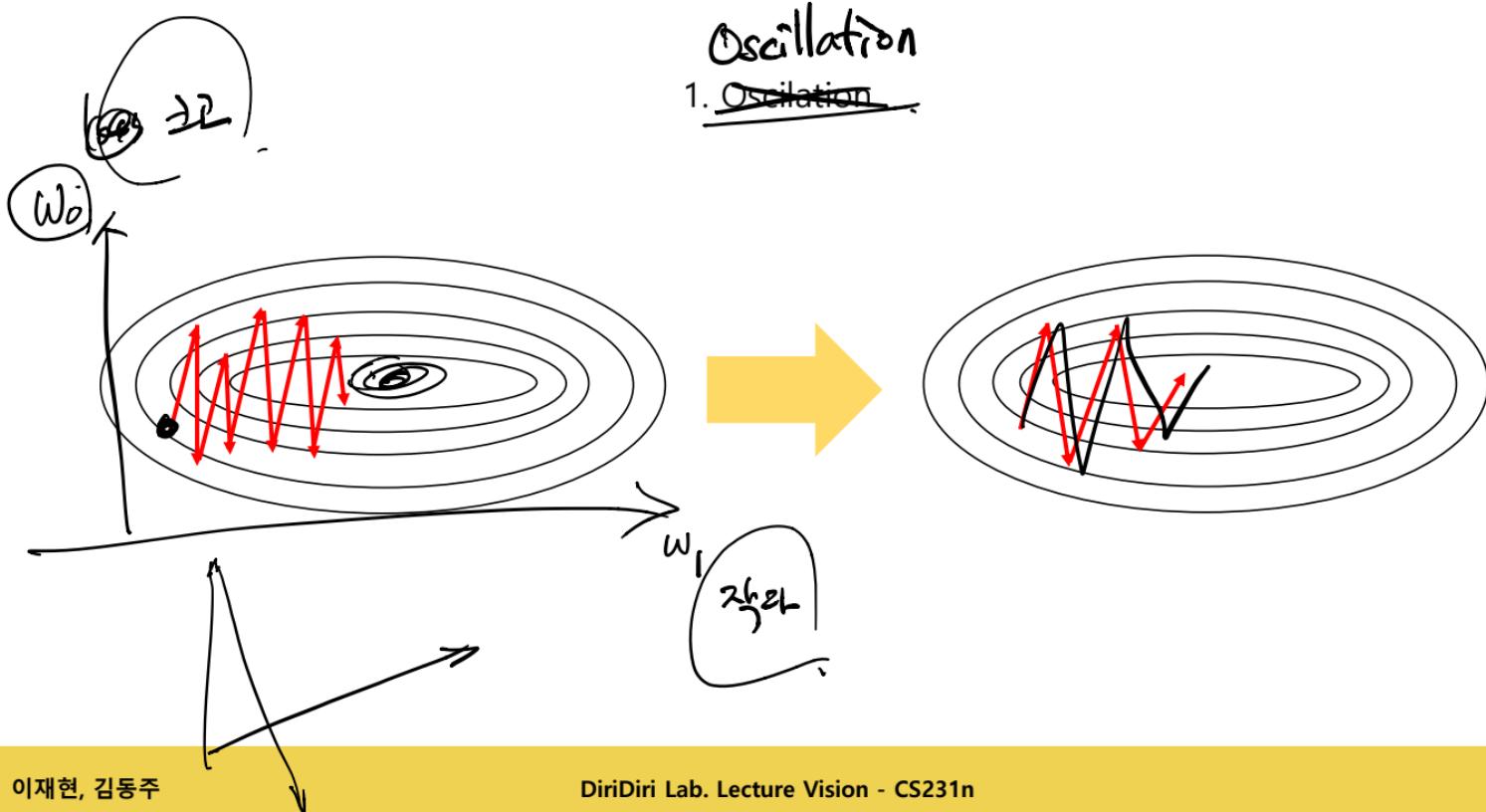


batch\_size = 32

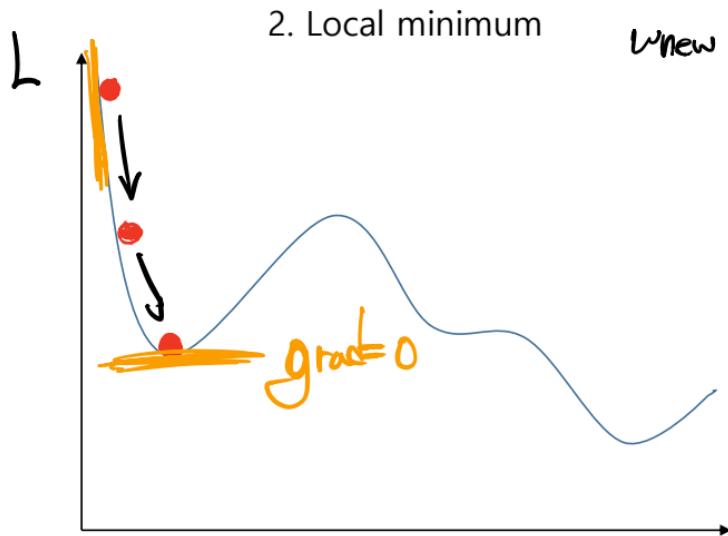
## Stochastic Gradient Descent (SGD)



## Two Problems of SGD



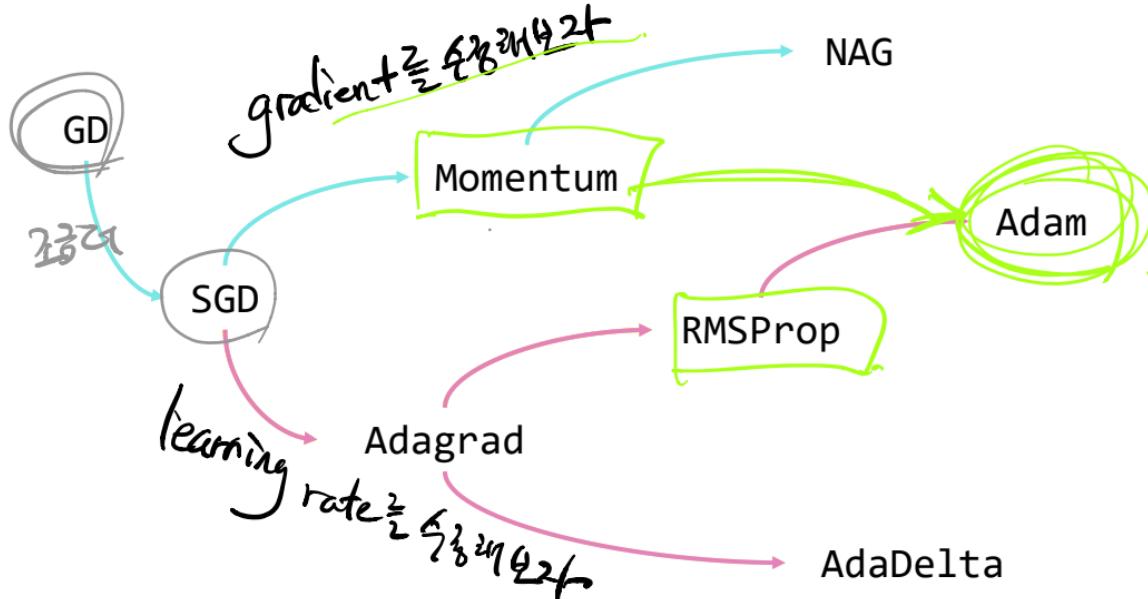
## Two Problems of SGD



2. Local minimum

$$\omega_{\text{new}} = \omega - \eta \times \cancel{\frac{\partial L}{\partial \omega}}$$
$$\omega_{\text{new}} = \omega_{\text{old}}$$

# Optimizer 발전 과정



# Momentum

Idea : gradient update에 관성을 더하자!

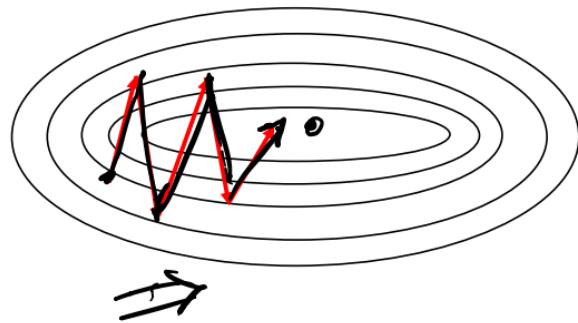
$$w_{new} = w - v_t \quad v_t = \text{속도 벡터}$$

$$v_t = \gamma v_{t-1} + \eta \frac{\partial L}{\partial w}$$

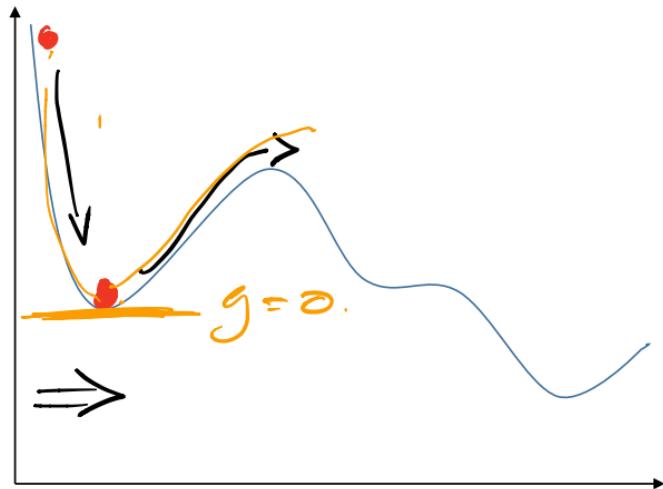
friction 이전의 속도  
마찰을 \_\_\_\_\_

$\gamma$ : 보통은 0.9

# Momentum



$$\mathbf{v}_f = \mathbf{r} \cdot \mathbf{\omega}_{\text{frame}} + \mathbf{\Omega} \times \mathbf{v}$$

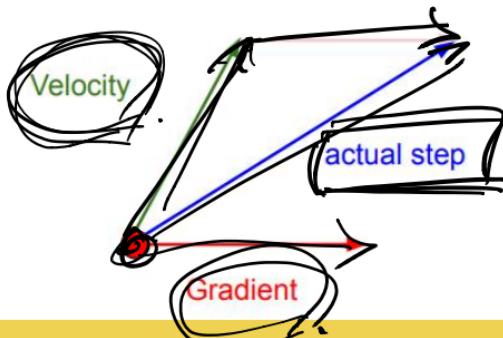


# Nestrov Accelerated Gradient (NAG)

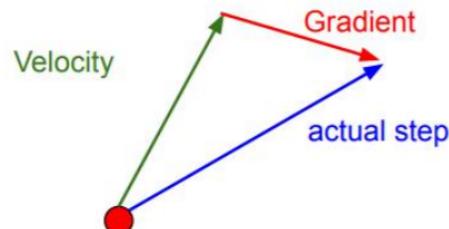
$$w_{new} = w - v_t$$

$$v_t = \gamma v_{t-1} + \eta \frac{\partial L}{\partial w} \Big|_{w=w-\gamma v_{t-1}}$$

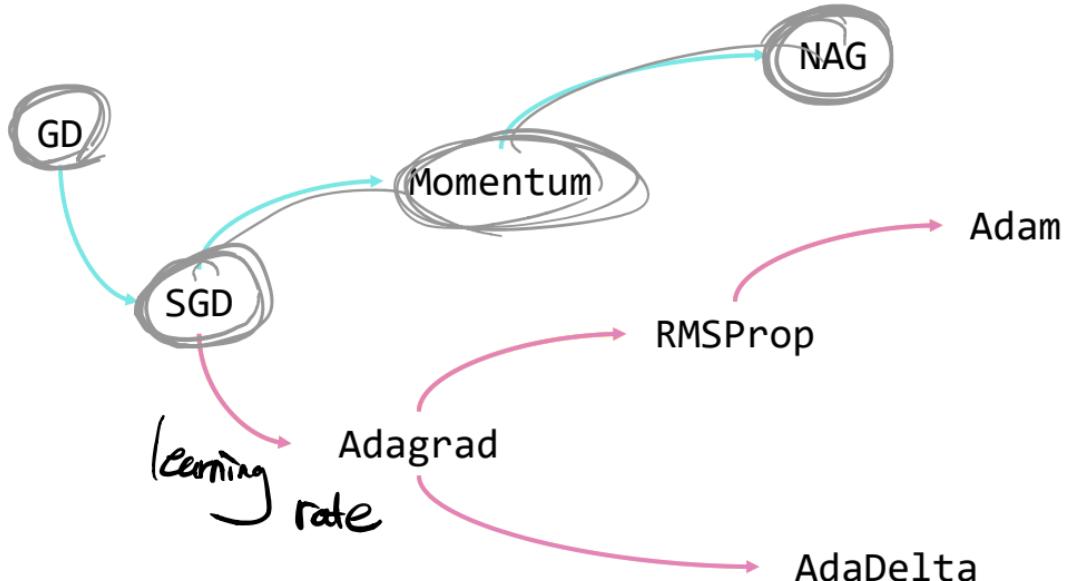
Momentum update:



Nesterov Momentum



# Optimizer 발전 과정



## Adagrad

$$\eta' = \frac{\eta}{\sqrt{G_t}}$$

Idea : learning rate를 gradient에 따라 다르게 탐색하자!

$w_0$	$w_1$
-5	0.1
-20	0.05

$$w_{new} = w - \frac{\eta}{\sqrt{G_t + \epsilon}} \frac{\partial L}{\partial w}$$

$$G_t = G_{t-1} + \left( \frac{\partial L}{\partial w} \right)^2$$

이전까지 gradient scale을 저장

$$G: \begin{bmatrix} 425 & 0.0125 \\ 20 & 0.1 \end{bmatrix}$$

$$w_0 \Rightarrow \frac{0.01}{20} = 0.0005$$

$$w_1 \Rightarrow \frac{0.01}{0.1} = 0.1$$

## Problem of Adagrad

$$G_t \rightarrow \infty \quad \eta' \rightarrow 0$$

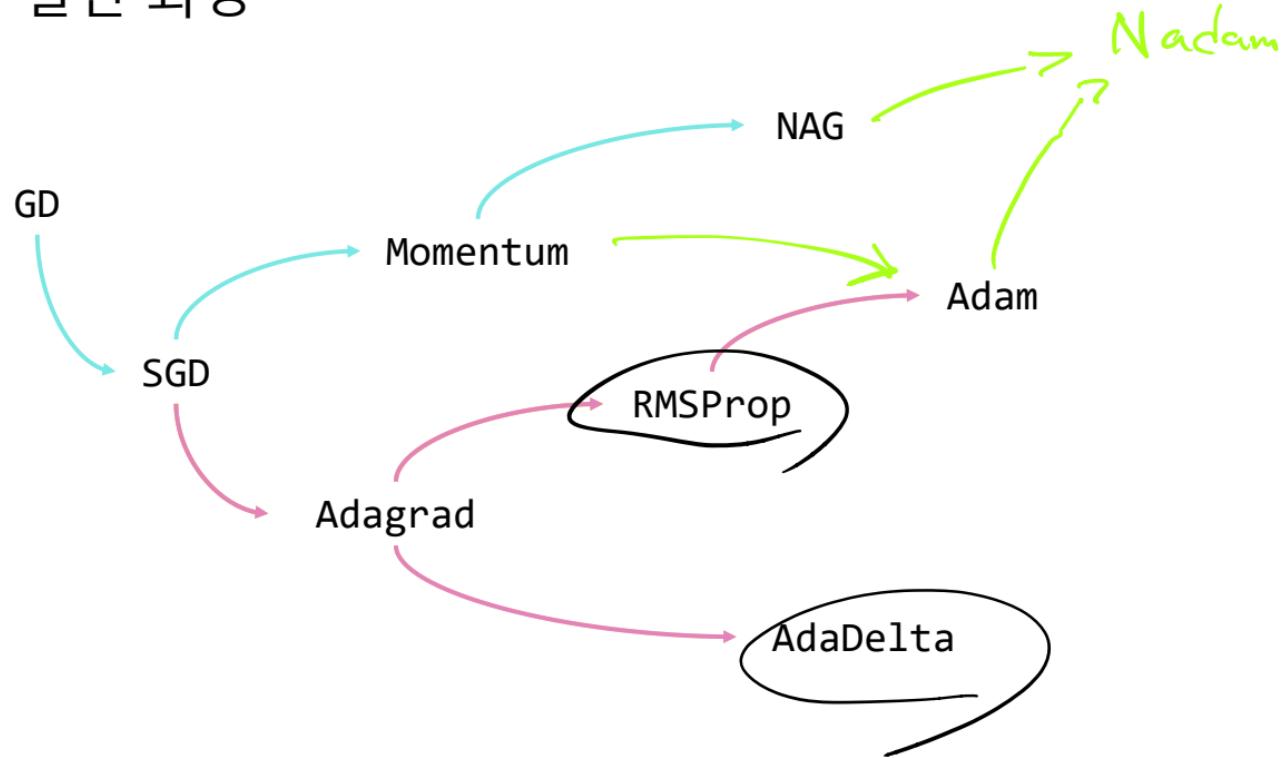
$$\eta' = \frac{\eta}{\sqrt{G_t + \epsilon}} \rightarrow \infty$$

$G_t$  가 무한정 커지면...?

→ learning rate  $\eta$ 가 점점 0에 가까워진다!

→ 학습이 안 된다..!

# Optimizer 발전 과정



## RMSProp

## Adagrad

Idea : gradient에 따라 learning rate 바꾸는 건 좋은데, 무한정 커지는 걸 방지하자!

$$G_t \leftarrow G_{t-1}$$

$$\frac{\partial L}{\partial w}^2$$

$$w_{new} = w - \frac{\eta}{\sqrt{G + \varepsilon}} \frac{\partial L}{\partial w}$$

$$G_t = \gamma G_{t-1} + (1 - \gamma) \left( \frac{\partial L}{\partial w} \right)^2$$

이동 EMA  
자수평균

$$\gamma : 0.99$$

## AdaDelta

$$\Delta w \propto \frac{\partial L}{\partial w} \propto \frac{1}{u(w)}$$

Idea : gradient에 따라 learning rate 바꾸는 건 좋은데, 무한정 커지는 걸 방지하자!

$$w_{new} = w - \Delta w$$

$$\Delta w = \frac{\sqrt{s + \varepsilon}}{\sqrt{G + \varepsilon}} \frac{\partial L}{\partial w}$$

$$\Delta w \propto \frac{\frac{\partial L}{\partial w}}{\sqrt{\frac{\partial L}{\partial w^2}}} \propto u(w)$$

$$G = \underbrace{\gamma G + (1 - \gamma) \left( \frac{\partial L}{\partial w} \right)^2}_{s = \gamma s + (1 - \gamma) \Delta w^2}$$

Adam  $\rightarrow$  Momentum + RMSProp.

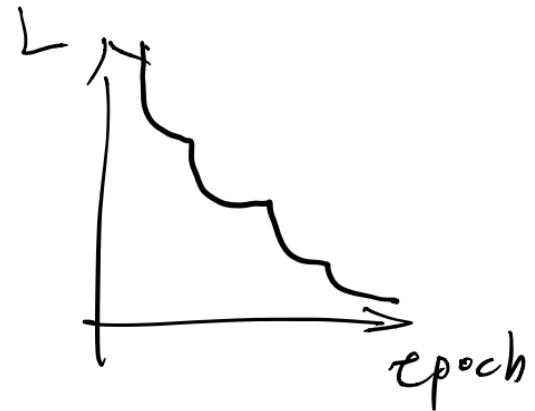
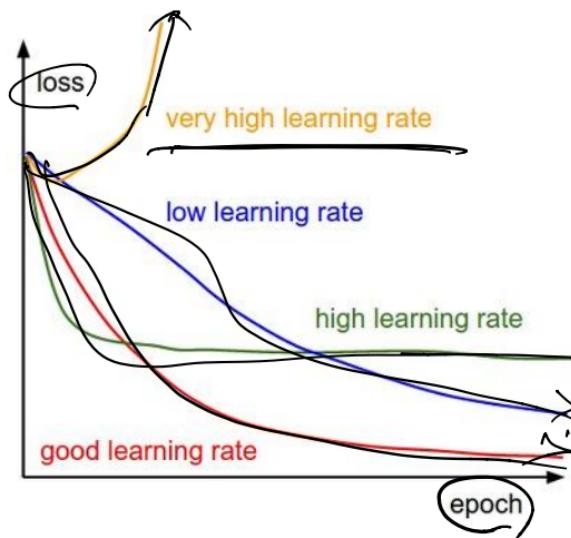
Idea : RMSProp에 Momentum을 합친다면..?

$$w_{new} = w - \frac{\eta}{\sqrt{\hat{v}_t + \epsilon}} \hat{m}_t \quad g = \frac{\partial L}{\partial w}$$
$$\hat{m}_t = \frac{m_t}{1 - \beta_1^t}, \quad \hat{v}_t = \frac{v_t}{1 - \beta_2^t}$$
$$\left( \begin{array}{l} m_t = \beta_1 m_{t-1} + (1 - \beta_1) \frac{\partial L}{\partial w} \\ v_t = \beta_2 v_{t-1} + (1 - \beta_2) \left( \frac{\partial L}{\partial w} \right)^2 \end{array} \right) \Rightarrow 0$$

## Learning rate decay

0.01  
0.1)

SGD + lr decay



# Momentum Nestrov ?

## Optimizers in Keras

tf. keras.optimizers.SGD(lr=0.01, momentum=0.0, decay=0.0, nesterov=False)

True  
SGD  
Momentum  
NAG

tf. keras.optimizers.RMSprop(lr=0.001, rho=0.9, epsilon=None, decay=0.0)

$\rho$        $\epsilon$

tf. keras.optimizers.Adagrad(lr=0.01, epsilon=None, decay=0.0)

$\epsilon$

tf. keras.optimizers.Adam(lr=0.001, beta\_1=0.9, beta\_2=0.999, epsilon=None, decay=0.0, amsgrad=False)

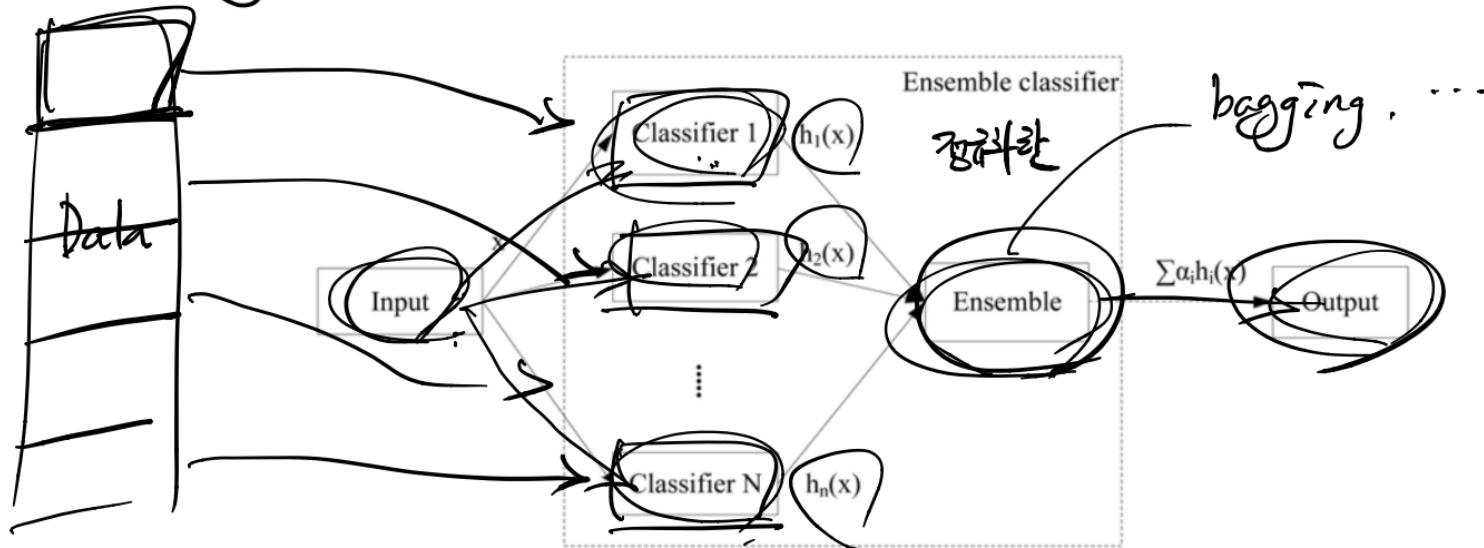
$\beta_1$        $\beta_2$        $\epsilon$

# Regularization

## Model Ensemble

overfitting

Idea : 모델 하나만 학습시킬 게 아니라, 여러 개를 학습시켜서 합하자!



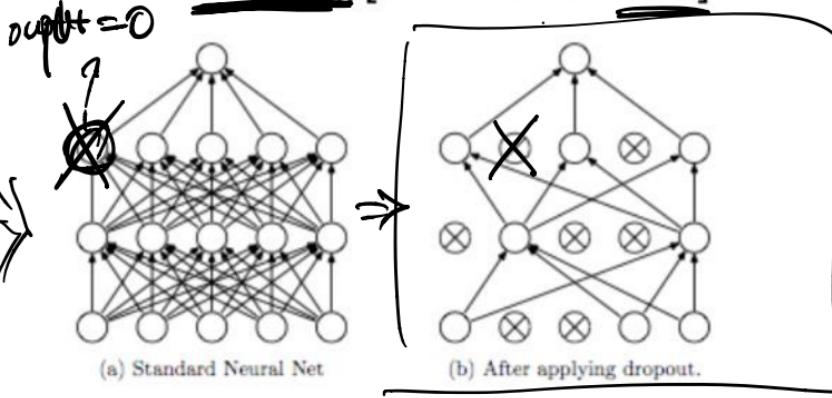
Dropout

$$\underline{p = 0.8}$$

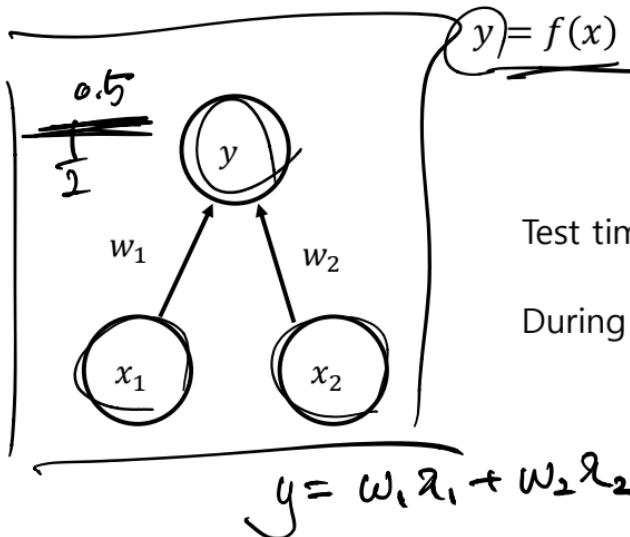
오늘  
암상

- → 끄기 있나?
- ✕ 털이 있나?
- → 눈이 놀라운가?
- :
- ✕

Dropout: A Simple Way to Prevent Neural Networks  
from Overfitting [Srivastava et al. 2014]



Dropout for test data..?



Test time:  $E[y] = \underline{w_1x_1 + w_2x_2}$

During Training:  $E[y] = \frac{1}{4}(w_1x_1 + w_2x_2) + \frac{1}{4}(w_1x_1 + 0) + \frac{1}{4}(0 + w_2x_2) + \frac{1}{4}(0 + 0)$

$\frac{1}{2}(w_1x_1 + w_2x_2)$

## Dropout in Keras

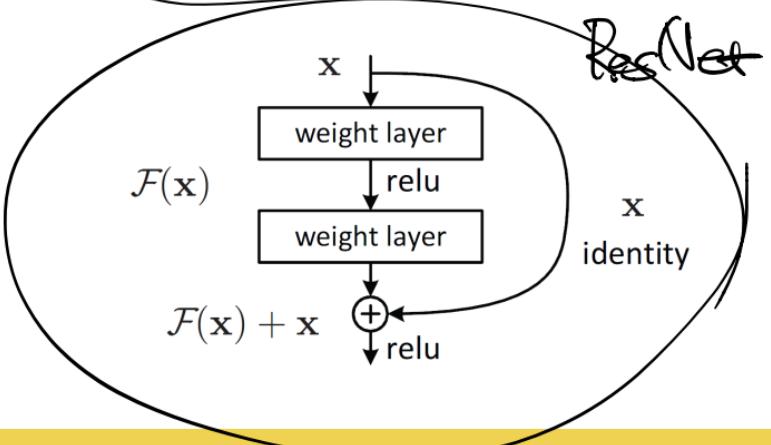
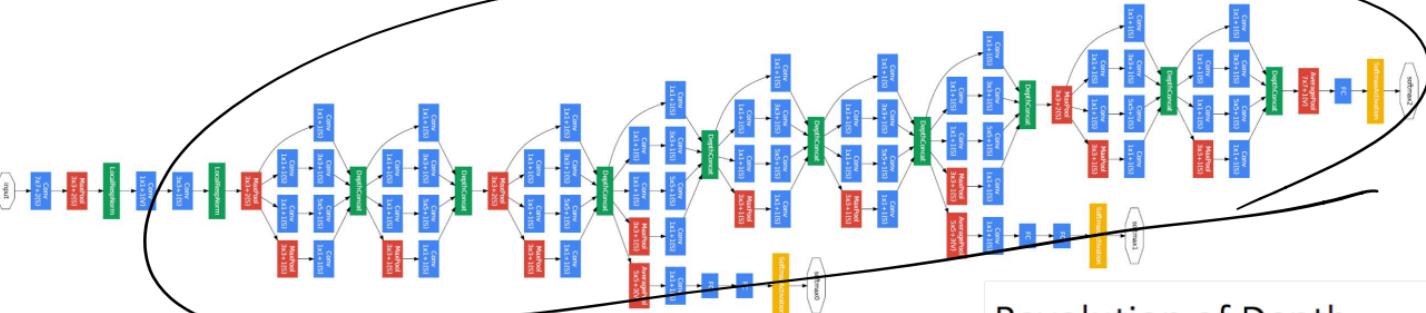
tf.keras.

model = Sequential([  
 tf.keras.layers.Flatten(),  
])  
 tf.keras.layers.Dropout(0.2)

tf.keras.layers.Dropout(rate, noise\_shape=None, seed=None)

# Preview on Next Class

GoogLeNet



## Revolution of Depth

152 layers

