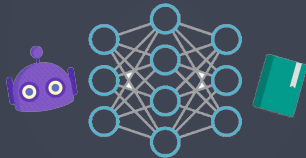


# Deep Learning

## Chapter 3 활성화 함수, 오차 역전파, 경사하강법 (Activation Function, Back Propagation, Gradient Descent Algorithm)



START

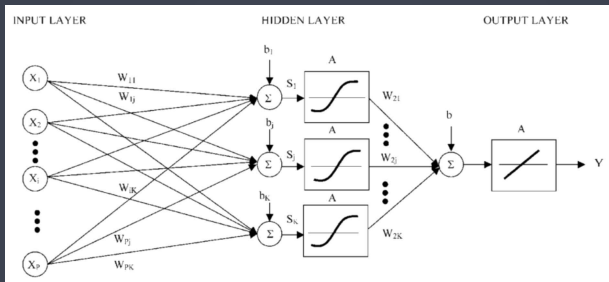


- 활성화 함수의 개념을 이해 하고 종류를 알 수 있다.
- 오차역전파의 개념을 이해 할 수 있다.
- 다양한 경사하강법 종류를 알 수 있다.
- Keras를 활용해 다양한 경사하강법을 적용 할 수 있다.

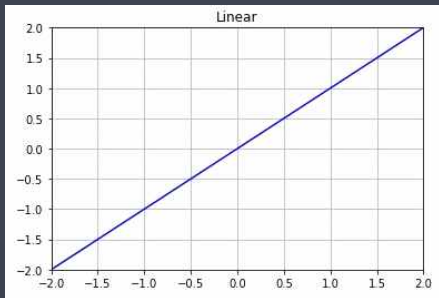
## 활성화 함수(Activation Function)

- 신경망은 한 계층의 신호를 다음 계층으로 그대로 전달하지 않고 활성화 함수를 거친 후에 전달함
- 사람의 신경망 속 뉴런들도 모든 자극을 다 다음 뉴런으로 전달하는 것은 아니고 **역치 이상의 자극만 전달**하게 됨
- 활성화 함수는 이런 부분까지 사람과 유사하게 구현하여 **사람처럼 사고하고 행동하는 인공지능 기술을 실현**하기 위해 도입됨
- 또한 선형모델을 기반으로 하는 딥러닝 신경망에서 **분류 문제를 해결**하기 위해서 비선형 활성화 함수가 필요함

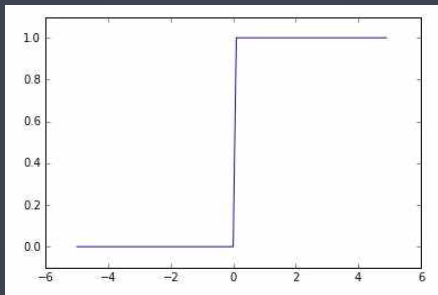
층에 따라 다른 활성화 함수를 사용 할 수 있다.



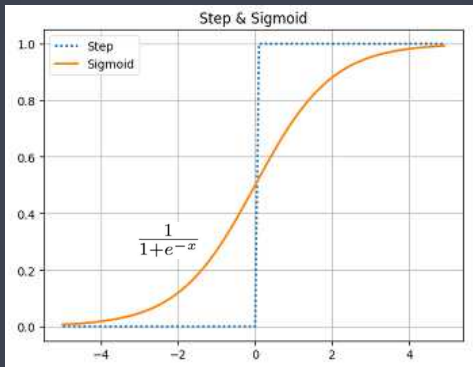
Linear function(선형함수=항등함수) → 회귀



Step function(계단 함수) → 분류의 초기 활성화 함수



## Sigmoid 함수 → 이진분류



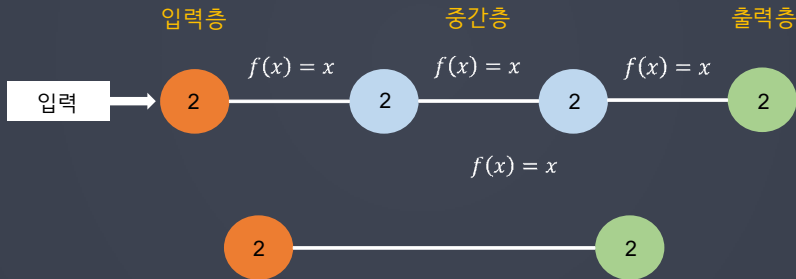
## 1. 중간층에 활성화 함수로 비선형 함수를 사용하는 이유

- 계단 함수(step)와 시그모이드 함수(sigmoid)는 비선형 함수이다.
- 중간층 활성화 함수로 선형함수(linear)를 사용하면 다층 구조의 효과를 살릴 수 없다.



## 2. 중간층에 활성화 함수로 선형 함수를 사용하게 된다면

선형함수(linear) 수식은  $h(x) = x$



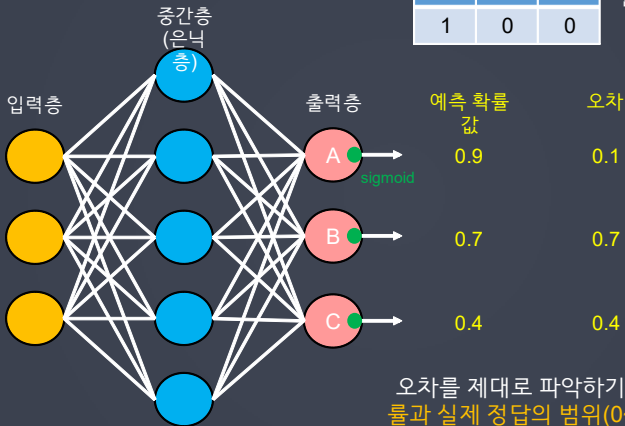
- 딥러닝 신경망에서 다중분류 문제를 해결하는 프로세스는 각 클래스에 대한 확률 값을 토대로 가장 높은 확률 값을 가지는 클래스로 최종 분류를 진행함
- 각 레이블의 확률들을 알기 위해 출력층 퍼셉트론 개수를 클래스 개수와 맞춰야 함(하나의 퍼셉트론이 하나의 클래스에 대한 확률 값을 출력)
- 또한 다중 분류 문제를 풀 경우 정답 데이터를 원 핫 인코딩 해야 함
- 신경망 학습을 위해서는 원 핫 인코딩 된 정보(0,1)와 출력층의 각 퍼셉트론의 예측한 확률(0, 1)과의 오차를 바탕으로 신경망의 활성화 함수를



# 활성화 함수(Activation Function) 정리

A	B	C
1	0	0

원래 인코딩된 정답 데이터  
(A가 정답이라고 가정)



오차를 제대로 파악하기 위해서는 예측 확률과 실제 정답의 범위(0~1)가 같아야 비교 가능



## 소프트맥스(softmax) 함수 → 다중분류

다중분류에서 레이블 값에 대한 각 퍼셉트론의 예측 확률의 합을 1로  
설정

sigmoid에 비해 예측 오차의 평균을 줄여주는 효과

$$y_k = \frac{\exp(a_k)}{\sum_{i=1}^n \exp(a_i)}$$

입력값들의  
지수함수의 합

## 소프트맥스(softmax) 함수 코드 구현

```
1 import numpy as np
2
3 def softmax(x):
4     e_x = np.exp(x-x.max())
5     return e_x/e_x.sum()
```

```
1 x = np.array([1.0,1.0,2.0])
2 x
```

```
array([1., 1., 2.])
```

```
1 y = softmax(x)
2 y
```

```
array([0.21194156, 0.21194156, 0.57611688])
```

```
1 y.sum()
```

```
1.0
```

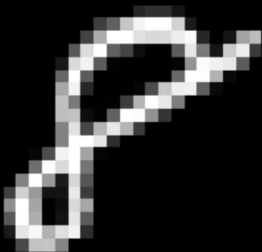
유 형	출력층 활성화 함수 (activation)	손실함수(=비용함수) (loss)
회귀	linear(항등 함수)	MSE
2진 분류	sigmoid(로지스틱 함수)	binary_crossentropy
다중 분류	softmax(소프트맥스 함수)	categorical_crossentropy

iris 데이터 신경망으로 풀기 (다중 분류)



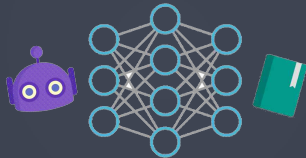
# Keras로 MNIST 손글씨 이미지 데이터 분류 모델을 만들어보자



[illegible]

# Keras로 패션 이미지 데이터 분류 모델을 만들어보자

# 오차 역전파 (Back Propagation)

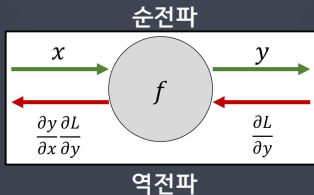


START



## 오차 역전파(Back Propagation)

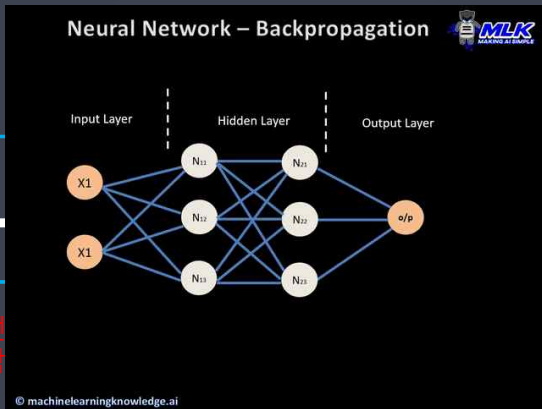
- **순전파** : 입력 데이터를 입력층에서부터 출력층까지 정 방향으로 이동시키며 출력 값을 **예측**해 나가는 과정
- **역전파** : 출력층에서 발생한 에러를 입력층 쪽으로 전파 시키면서 최적의 결과를 **학습**해 나가는 과정



MLP,

Data  
(입력데이터)

역전  
= 학

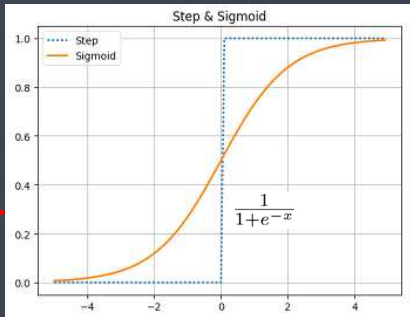
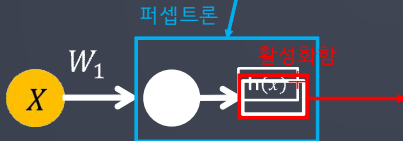


performance  
measure

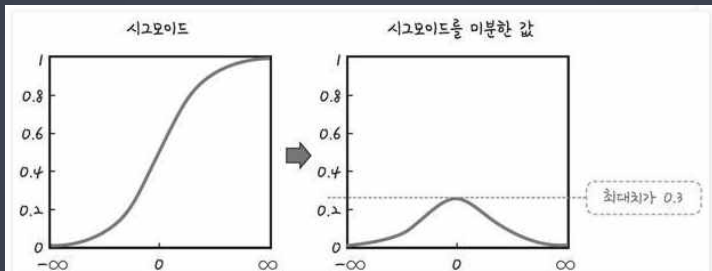
## 손실함수 및 Sigmoid 함수의 미분

- 신경망이 학습하기 위해서는 경사하강법(loss 함수를 미분)을 사

$$\text{용. } MSE = \frac{1}{m} \sum_{i=1}^m \text{예측값} (H(x_i) - y_i)^2$$

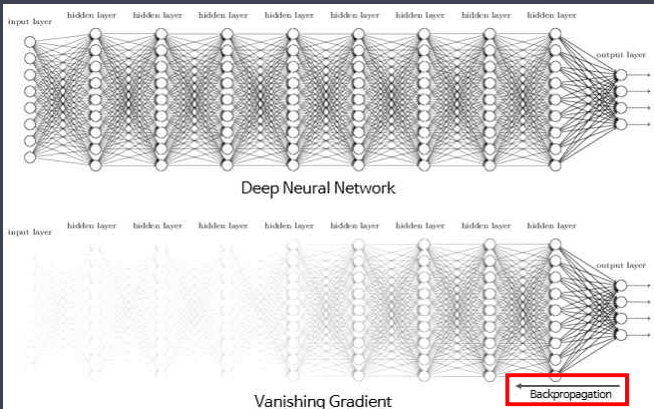


## Sigmoid 함수의 문제점



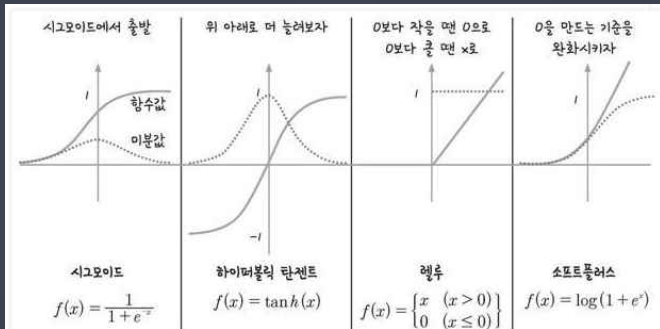
# Sigmoid 함수의 문제점

## - 기울기 소실 문제(Vanishing Gradient)



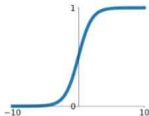


# 활성화 함수(Activation)의 종류

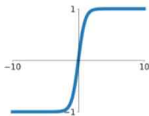


**Sigmoid**

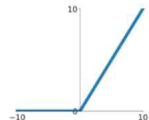
$$\sigma(x) = \frac{1}{1+e^{-x}}$$

**tanh**

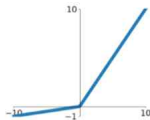
$$\tanh(x)$$

**ReLU**

$$\max(0, x)$$

**Leaky ReLU**

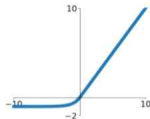
$$\max(0.1x, x)$$

**Maxout**

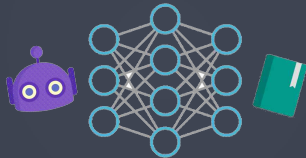
$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

**ELU**

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



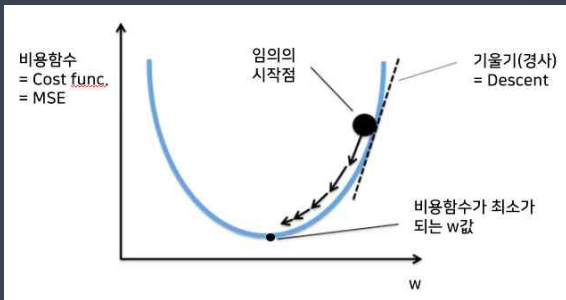
# 최적화 함수 (Optimizer)



START



## 경사하강법(Gradient Descent Algorithm)



# 최적화함수(Optimizer)의 종류



경사하강법  
(Gradient Descent)

전체 데이터를 이용해 업데이트

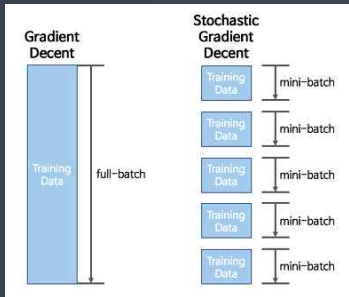


확률적경사하강법  
(Stochastic Gradient Descent)

확률적으로 선택된 일부 데이터를  
이용해 업데이트

## Batch\_size

일반적으로 PC 메모리의 한계 및 속도 저하 때문에 대부분의 경우에는  
한번의 epoch에 모든 데이터를 한꺼번에 집어넣기가 힘들



- batch\_size를 줄임
  - 메모리 소모가 적음(저 사양일 경우), 학습 속도가 느림, 정확도↑
- batch\_size를 높임
  - 메모리 소모가 큼, 학습 속도가 빠름, 정확도↓

→ batch\_size의 디폴트 값은 32이며  
일반적으로 32, 64가 많이 사용됨

# 최적화함수(Optimizer)의 종류



확률적경사하강법  
(Stochastic Gradient Descent)

확률적으로 선택된 일부 데이터를  
이용해 업데이트



모멘텀  
(Momentum)

경사 하강법에 관성을 적용해 업데이  
트

현재 batch뿐만 아니라 이전 batch  
데이터의 학습 결과도 반영

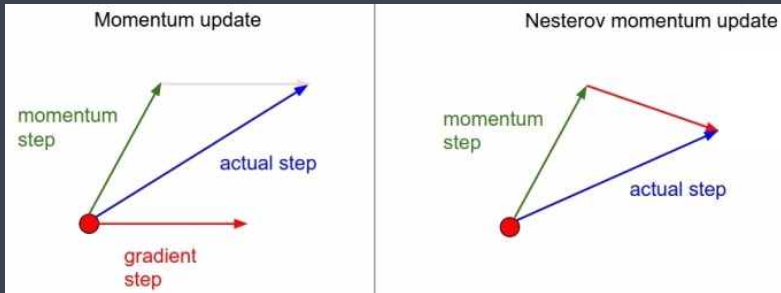
## 특징 (Momentum)

- 가중치를 수정하기 전 **이전 방향을 참고**하여 업데이트
- 지그재그 형태로 이동하는 현상이 줄어든다
- $\alpha$ 는 Learning Rate,  $m$ 은 momentum 계수 (보통 0.9)



$$V(t) = m * V(t - 1) - \alpha \frac{\partial}{\partial w} Cost(w)$$
$$W(t + 1) = W(t) + V(t)$$





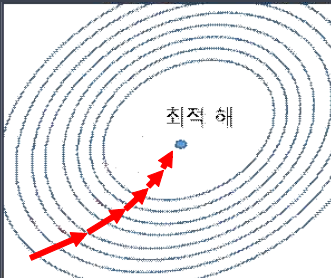
네스테로프 모멘텀  
(Nesterov Accelerated Gradient)  
개선된 모멘텀 방식

## 특징 (NAG)

- w, b값 업데이트 시 모멘텀 방식으로 먼저 더한 다음 계산
- 미리 해당 방향으로 이동한다고 가정하고 기울기를 계산해본 뒤 실제 업데이트 반영

- 분할수렴이든 아니든 무조건 수렴한다

$$V(t) = m * V(t - 1) - \alpha \frac{\partial}{\partial (w + m * V(t - 1))} Cost(w)$$
$$W(t + 1) = W(t) + V(t)$$

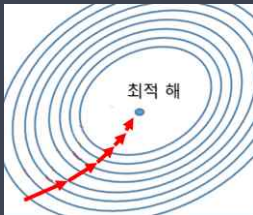


에이다그래드  
(Adaptive Gradient)

학습률 감소 방법을 적용해 업데이트

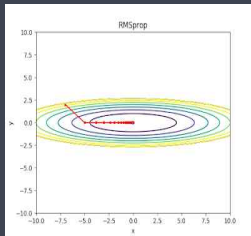
## 특징 (Adagrad)

- 학습을 진행하면서 **학습률을 점차 줄여가는** 방법
- 처음에는 크게 학습하다가 조금씩 작게 학습한다
- 학습을 빠르고 정확하게 할 수 있다



$$G(t) = G(t-1) + \left( \frac{\partial}{\partial w(t)} \text{Cost}(w(t)) \right)^2$$
$$= \sum_{i=0}^t \left( \frac{\partial}{\partial w(i)} \text{Cost}(w(i)) \right)^2$$

$$W(t+1) = W(t) - \alpha * \frac{1}{\sqrt{G(t)+\epsilon}} * \frac{\partial}{\partial w(i)} \text{Cost}(w(i))$$

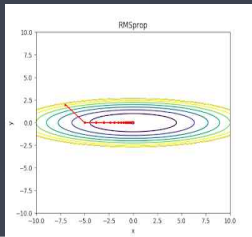


## RMSProp

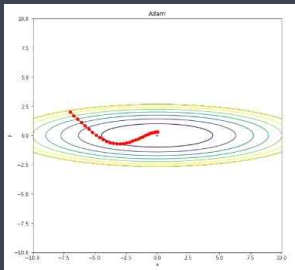
Adagrad의 단점을 해결한 최적화함수

## 특징 (RMSProp)

- Adagrad와 동일하게 학습을 진행하면서  
학습률을 점차 줄여가는 방법
- 최소값을 찾기전 학습이 멈추는 Adagrad의 단점을  
지수이동 평균을 도입해서 해결
- 지수 이동 평균 : 최근 학습한 수치의 영향력은 높이고  
과거 학습한 수치의 영향력은 낮추는 방식.



$$h_i \leftarrow \rho h_{i-1} + (1 - \rho) \frac{\partial L_i}{\partial W} \odot \frac{\partial L_i}{\partial W}$$

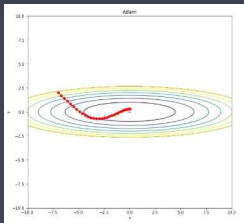


# Adam

Momentum과 RMSProp의 장점만을  
취한 최적화 함수

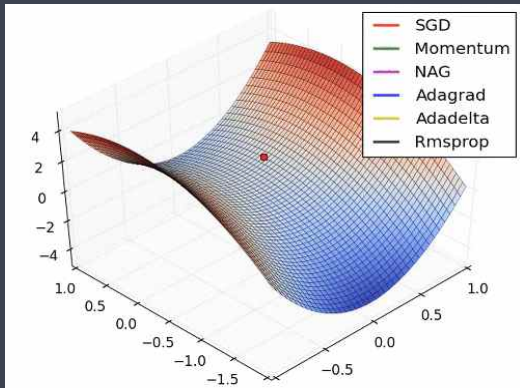
## 특징 (Adam)

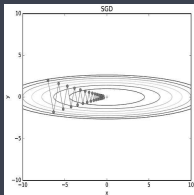
- 관성 방향으로 움직이는 Momentum과 보폭을 조절하며 움직이는 RMSProp의 장점을 하나로 합친 최적화 함수.
- 현재 보편적으로 사용하는 최적화 함수이며, 성능적인 측면에서 가장 나은 최적화 함수라 할 수 있다.



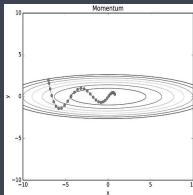
$$\begin{aligned} m_1 &\leftarrow \beta_1 m_0 + (1 - \beta_1) g_1 \\ \widehat{m}_1 &\leftarrow \frac{m_1}{1 - \beta_1^1} = \frac{\beta_1 m_0}{1 - \beta_1^1} + \frac{(1 - \beta_1) g_1}{1 - \beta_1^1} \\ &= 0 + g_1 (\because m_0 = 0) \end{aligned}$$



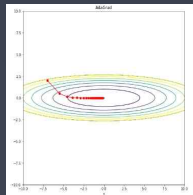




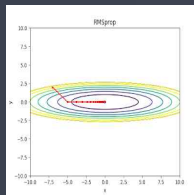
SGD



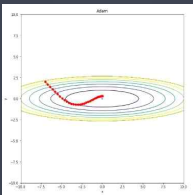
Momentum



Adagrad



RMSProp



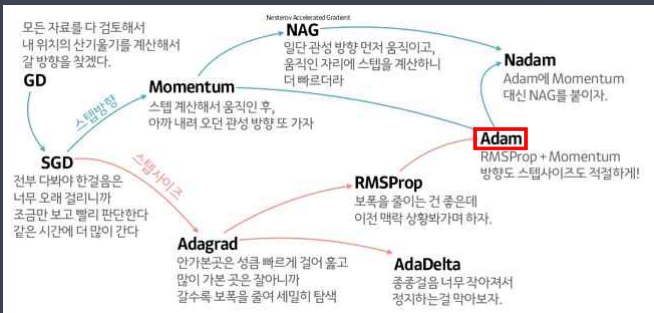
Adam

모든 문제에서 뛰어난 기법은 없다.  
하지만 많은 사람들이 Adam에  
만족하며 사용하는 것 같다.

But! 각자의 상황을 고려해보고,  
여러가지를 시도하는 것이  
가장 좋은 방법론이라고 말할 수 있다.

하지만 모르겠으면 Adam을 쓰는 것도 하나의  
방법론!

# 최적화함수(Optimizer)의 종류



출처 : <https://www.slideshare.net/yongho/ss-79607172>



## Keras

```
from tensorflow.keras import optimizers  
  
opti = optimizers.SGD(learning_rate=0.01, momentum=0.9)  
  
model.compile(loss='mse', optimizer=opti, metrics=['acc'])
```

Momentum

```
from tensorflow.keras import optimizers  
  
opti = optimizers.SGD(learning_rate=0.01, momentum=0.9, nesterov=True)  
  
model.compile(loss='mse', optimizer=opti, metrics=['acc'])
```

NAG

```
model.compile(loss="mse", optimizer="Adam", metrics=["acc"])
```

Adam

Adagrad, RMSprop, Adam 등은 이름으로 지정 가능



활성화함수, 최적화함수를 바꿔가며  
패션 이미지 데이터 분류 모델을 만들어  
비교해보자