

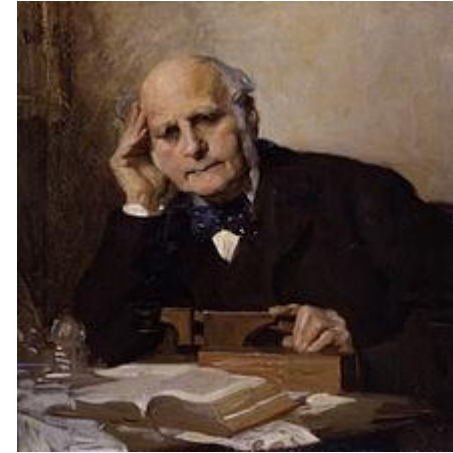
구현으로 배우는 딥러닝 이해

강사명 : 김정인

email : jikim@imguru.co.kr

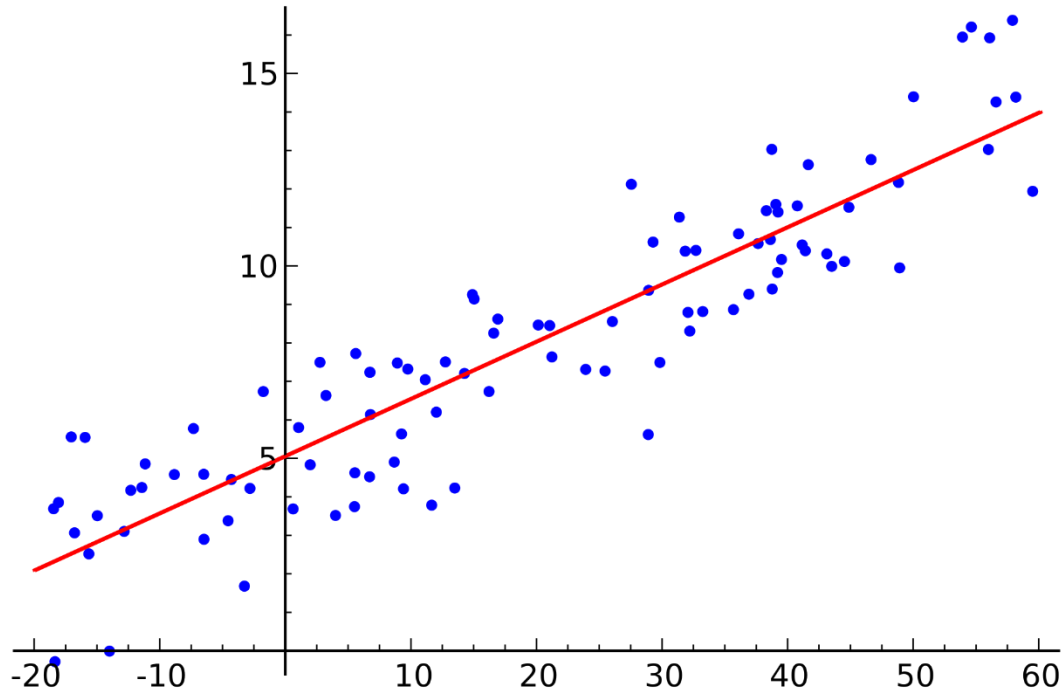
Regression - 회귀

"Regression toward the mean"



Sir Francis Galton
(1822 ~ 1911)

Linear Regression - 선형 회귀



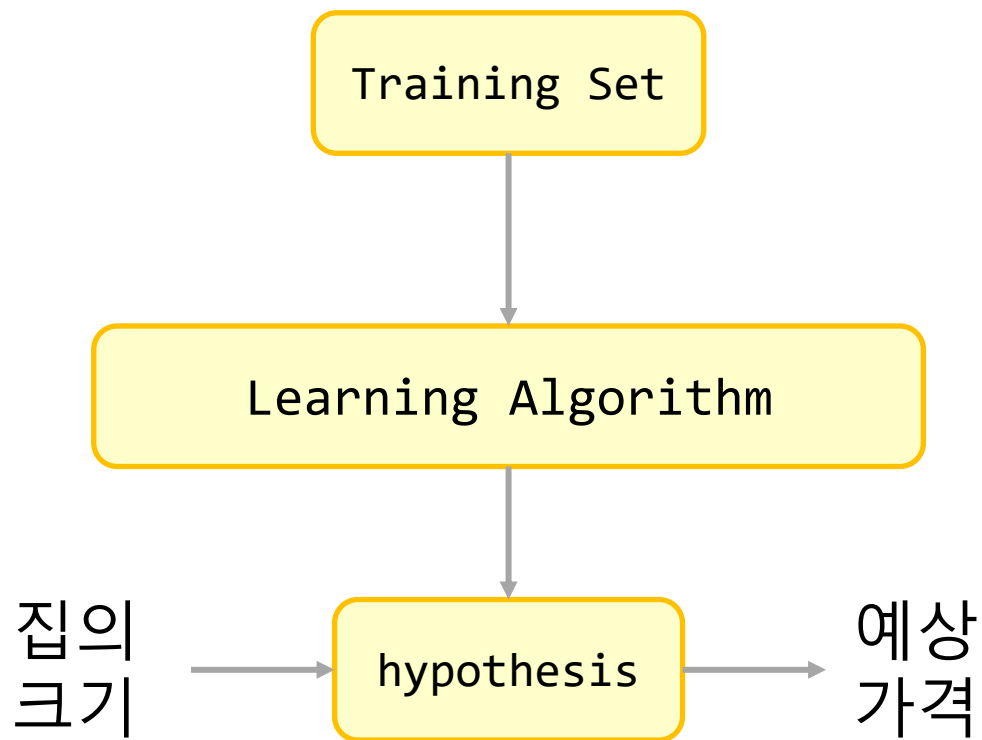
$$y = ax + b$$

↓

$$y = wx + b$$

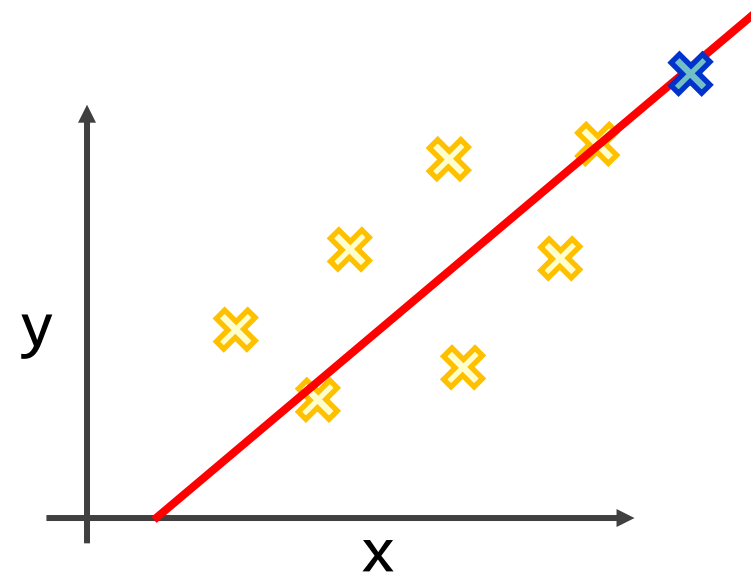
https://en.wikipedia.org/wiki/Linear_regression

Cost function



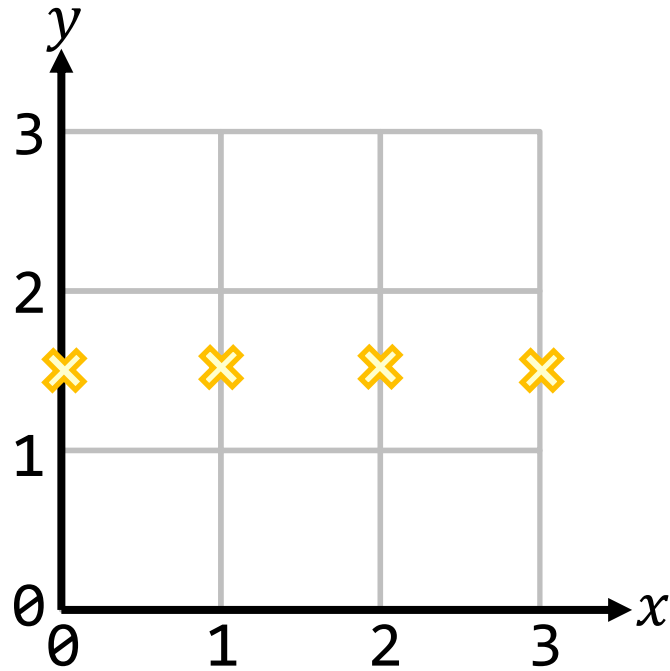
hypothesis ?

$$\hat{y} = wx + b$$

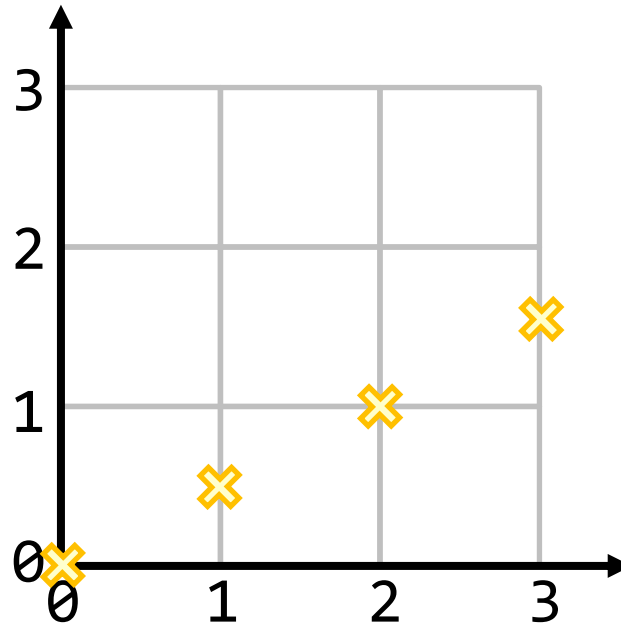


Cost function

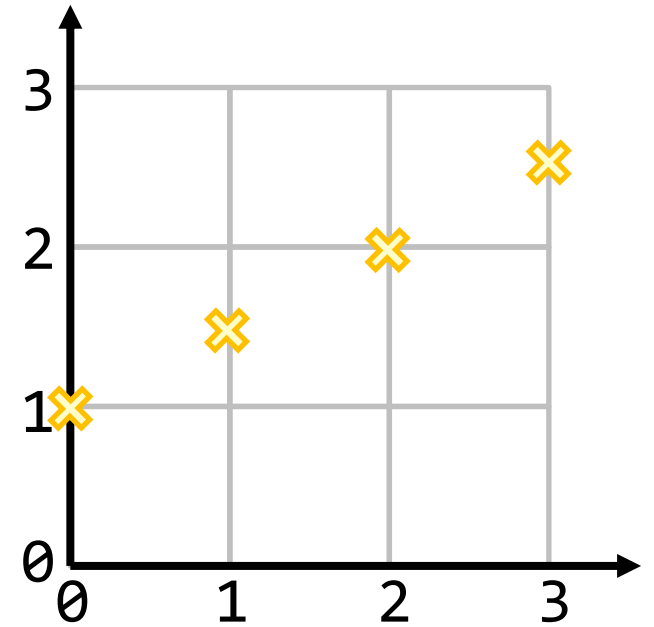
$$\hat{y} = wx + b$$



$$w = 0$$
$$b = 1.5$$



$$w = 0.5$$
$$b = 0$$



$$w = 0.5$$
$$b = 1$$

Hypothesis :

$$\hat{y} = wx + b$$

Parameters:

$$w, b$$

Cost Function

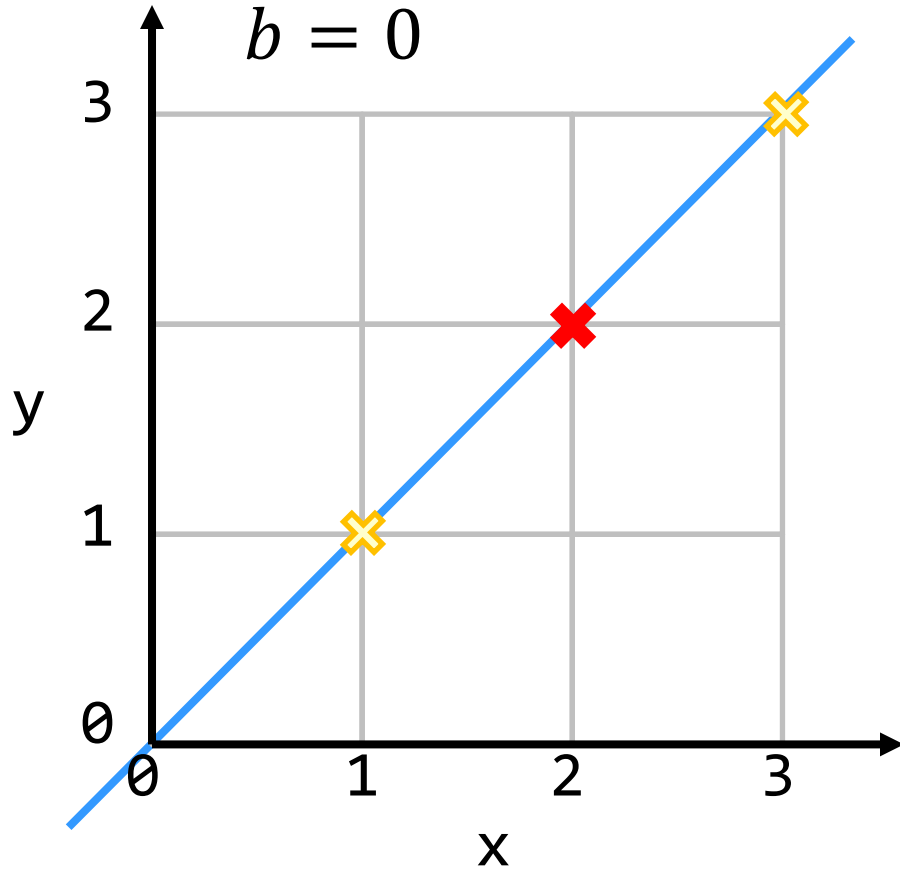
$$J(w, b) = \frac{1}{2} (\hat{y} - y)^2$$

Cost function

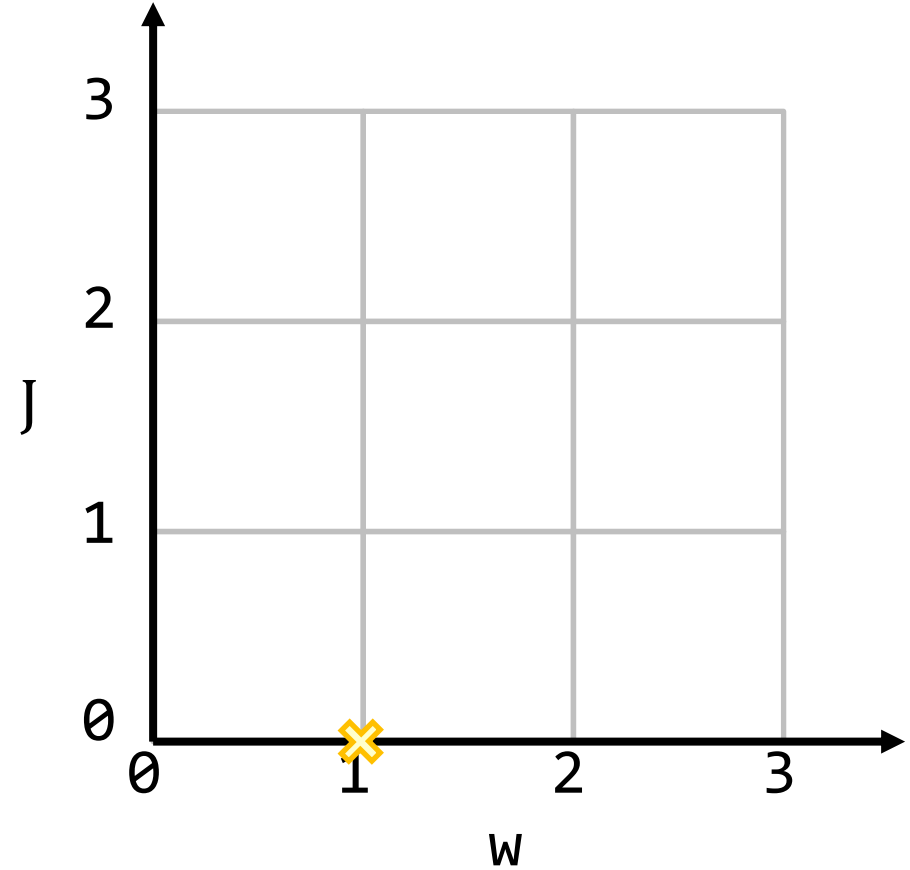
$$\hat{y} = wx + b$$

$$w = 1$$

$$b = 0$$



$$J(w, b) = \frac{1}{2} (\hat{y} - y)^2$$

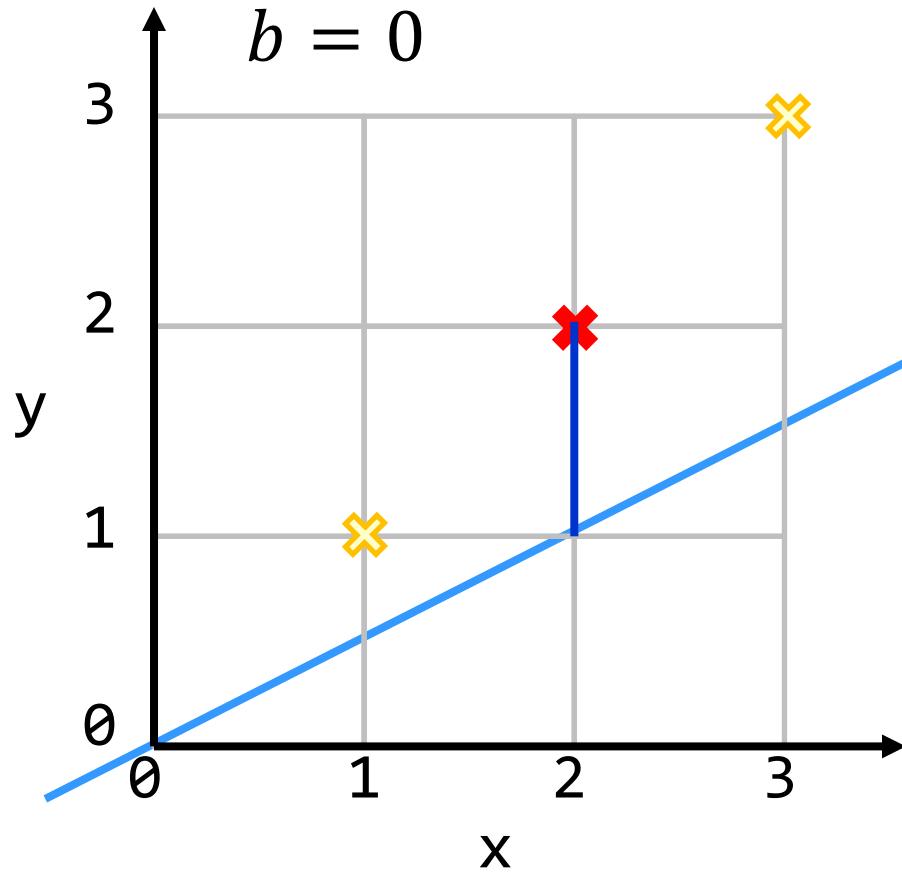


Cost function

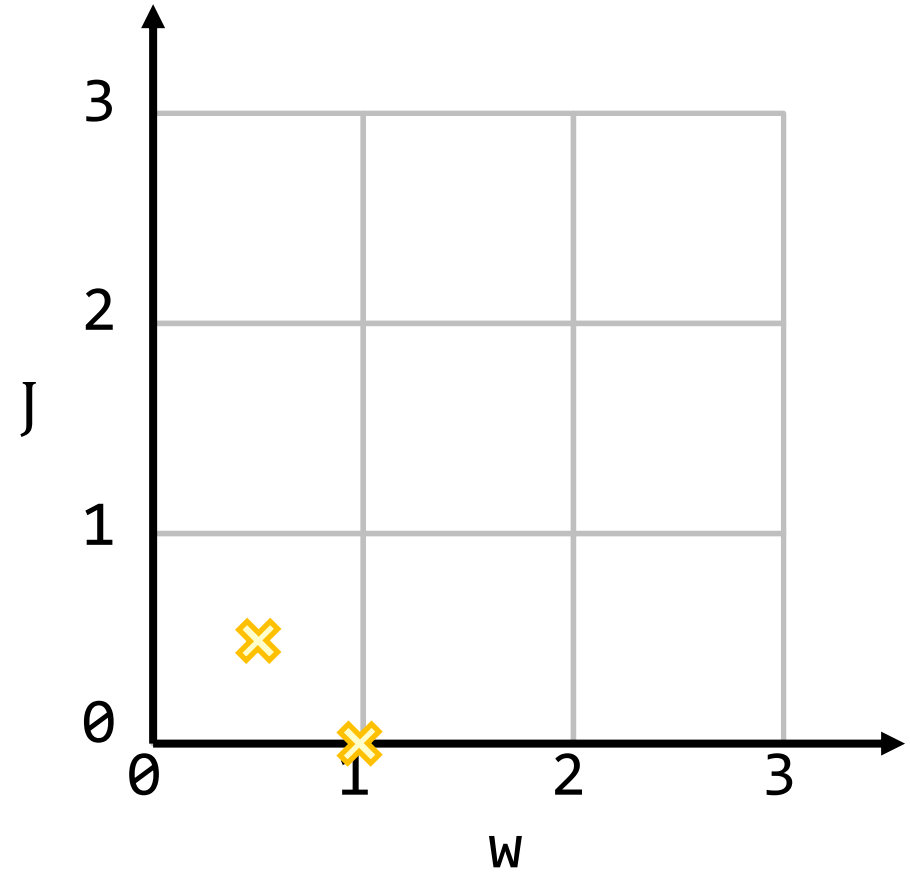
$$\hat{y} = wx + b$$

$$w = 0.5$$

$$b = 0$$



$$J(w, b) = \frac{1}{2} (\hat{y} - y)^2$$

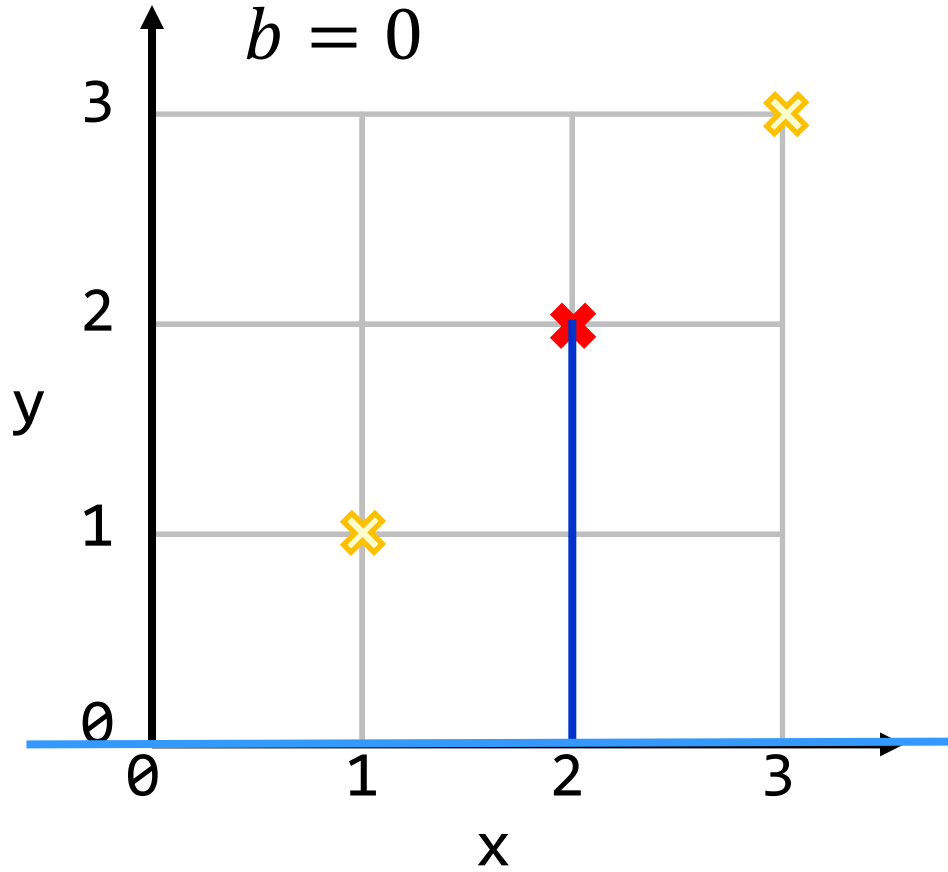


Cost function

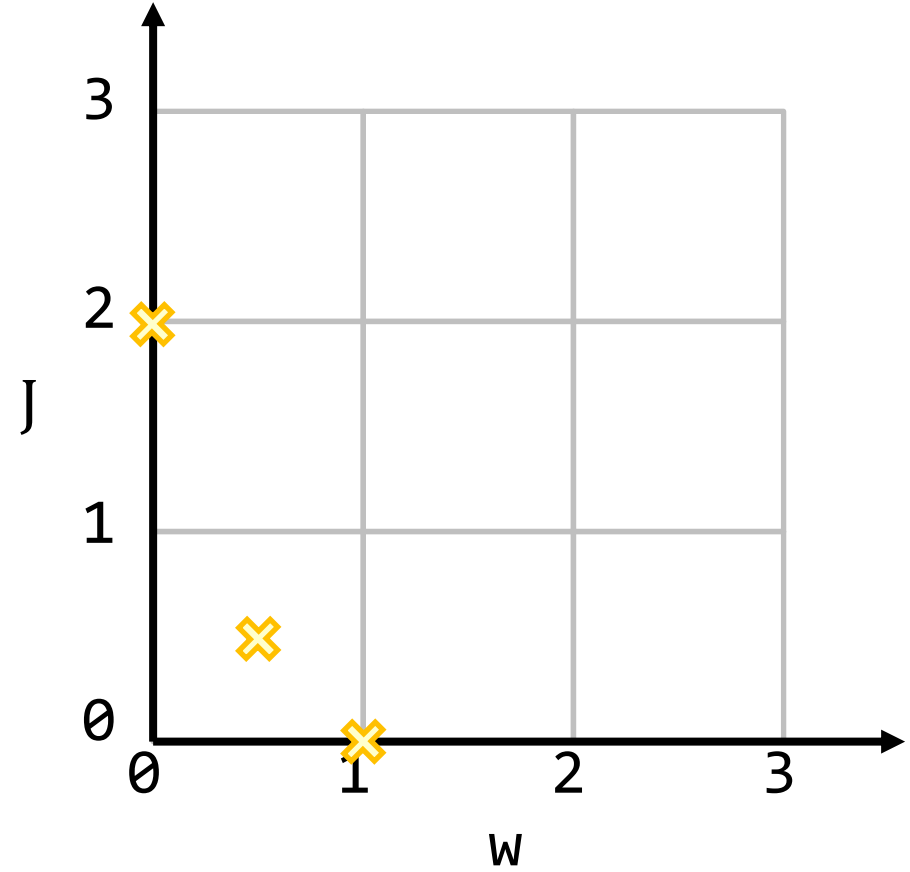
$$\hat{y} = wx + b$$

$$w = 0$$

$$b = 0$$



$$J(w, b) = \frac{1}{2} (\hat{y} - y)^2$$

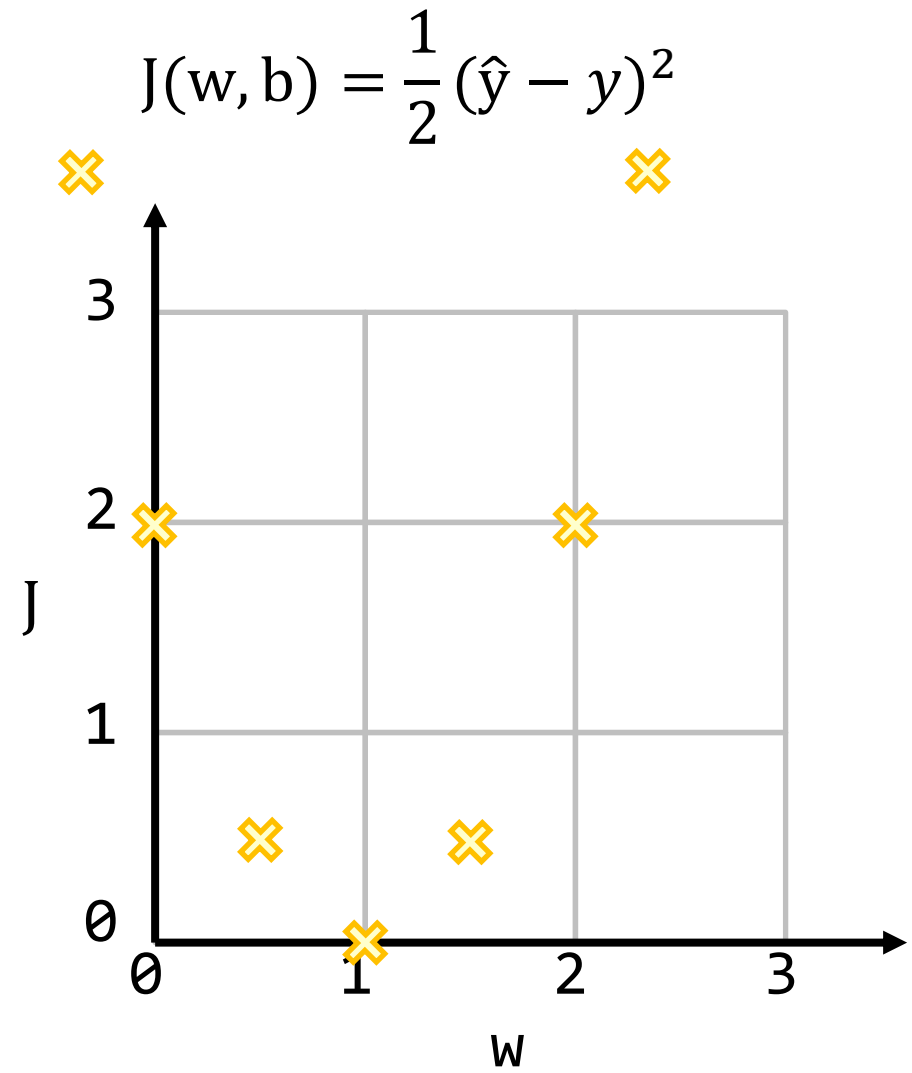
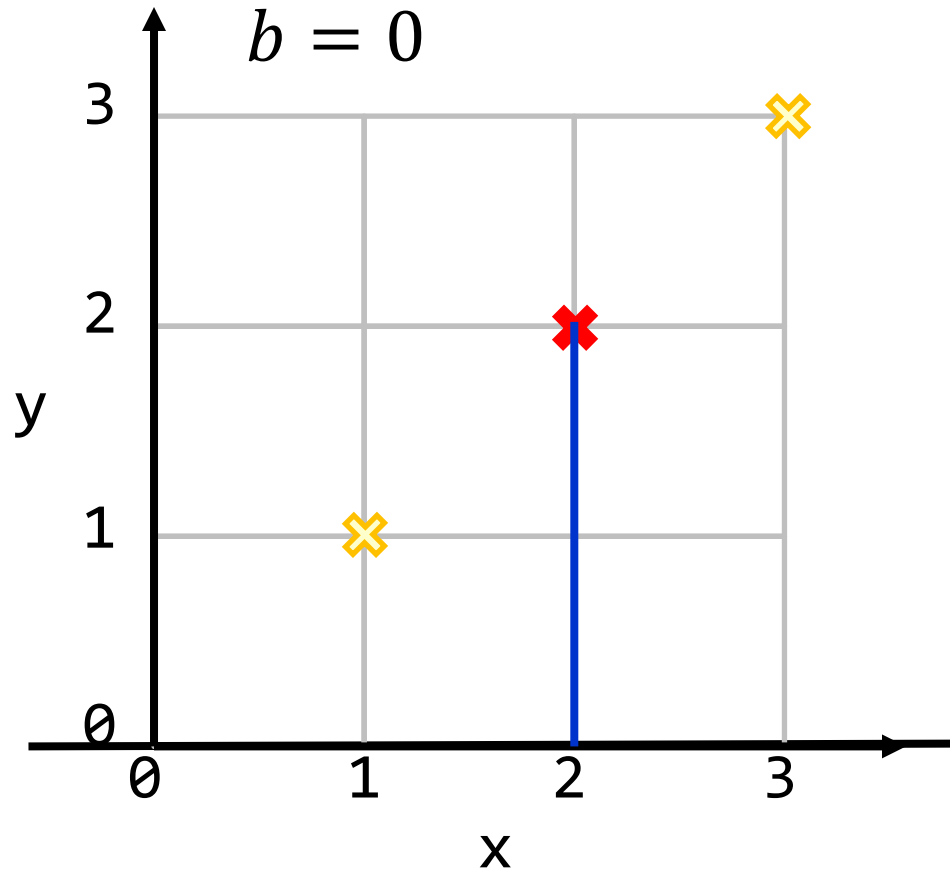


Cost function

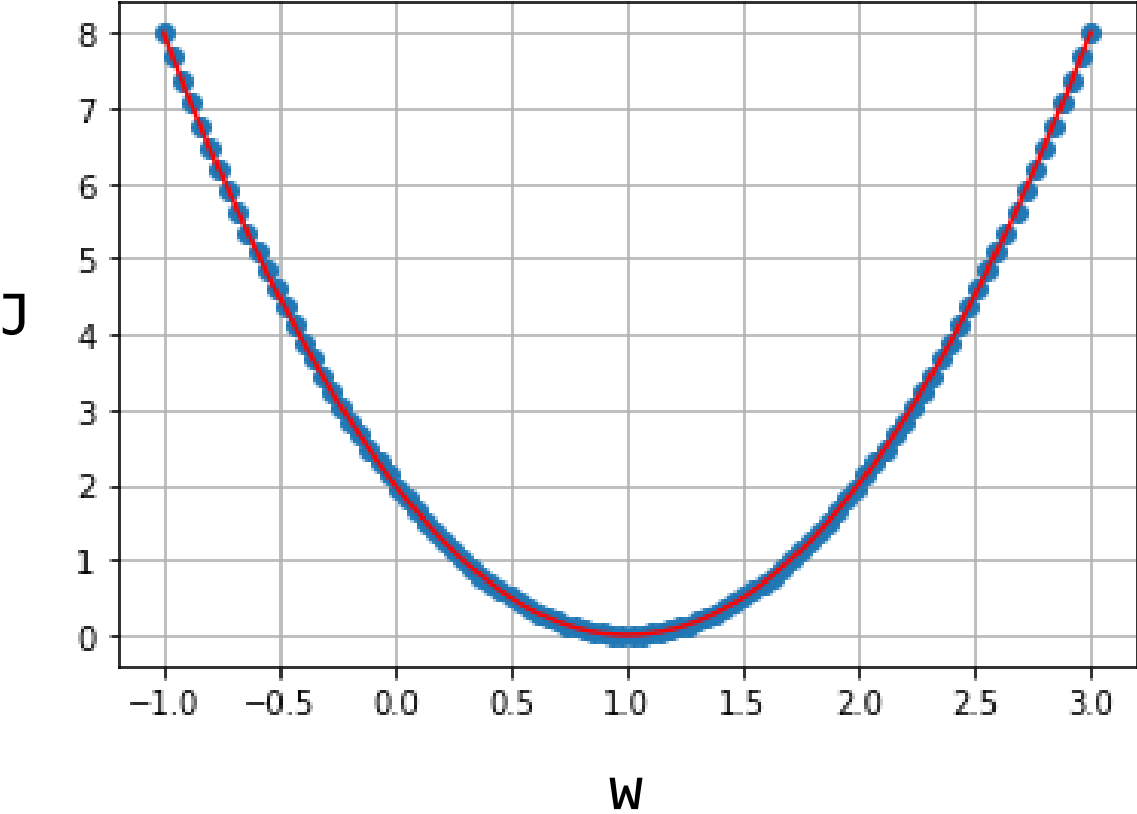
$$\hat{y} = wx + b$$

$$w = 0$$

$$b = 0$$



Cost function

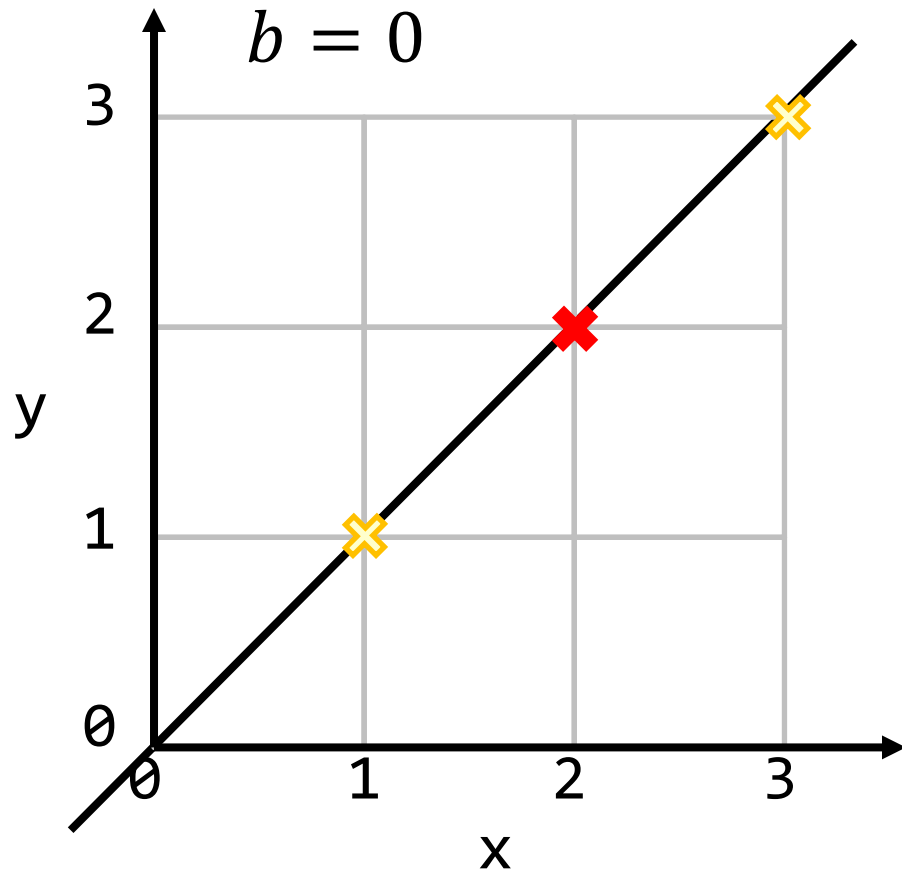


Cost function

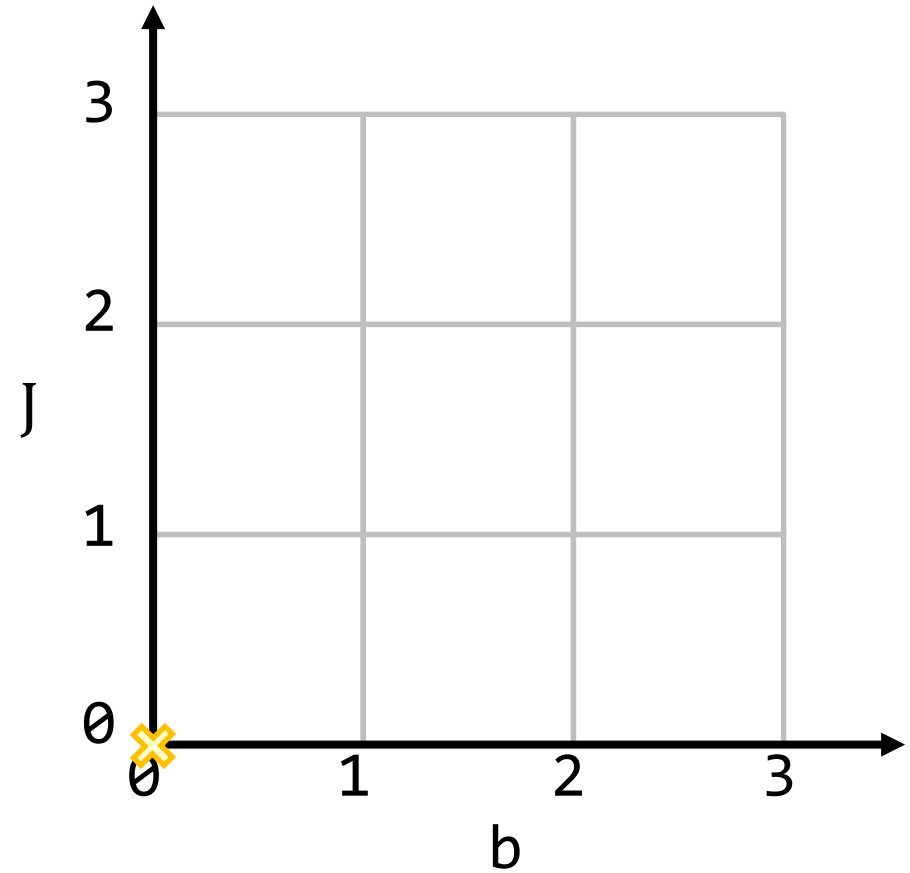
$$\hat{y} = wx + b$$

$$w = 1$$

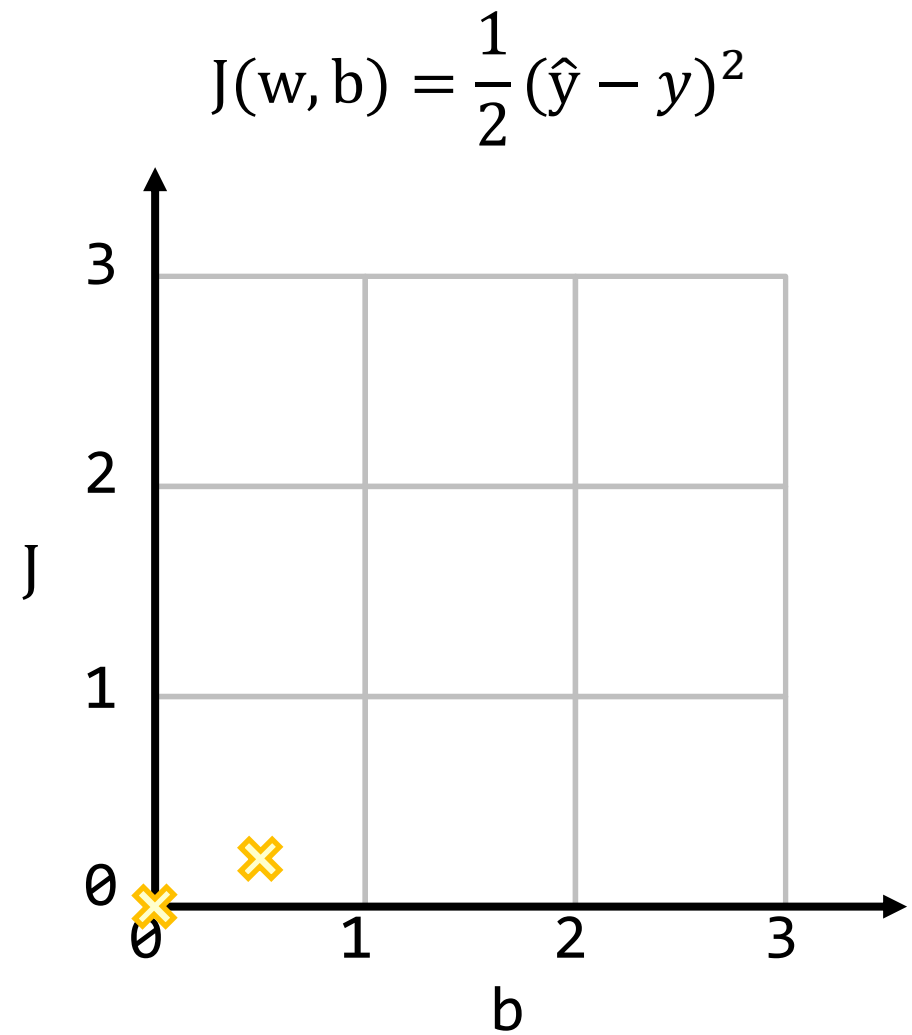
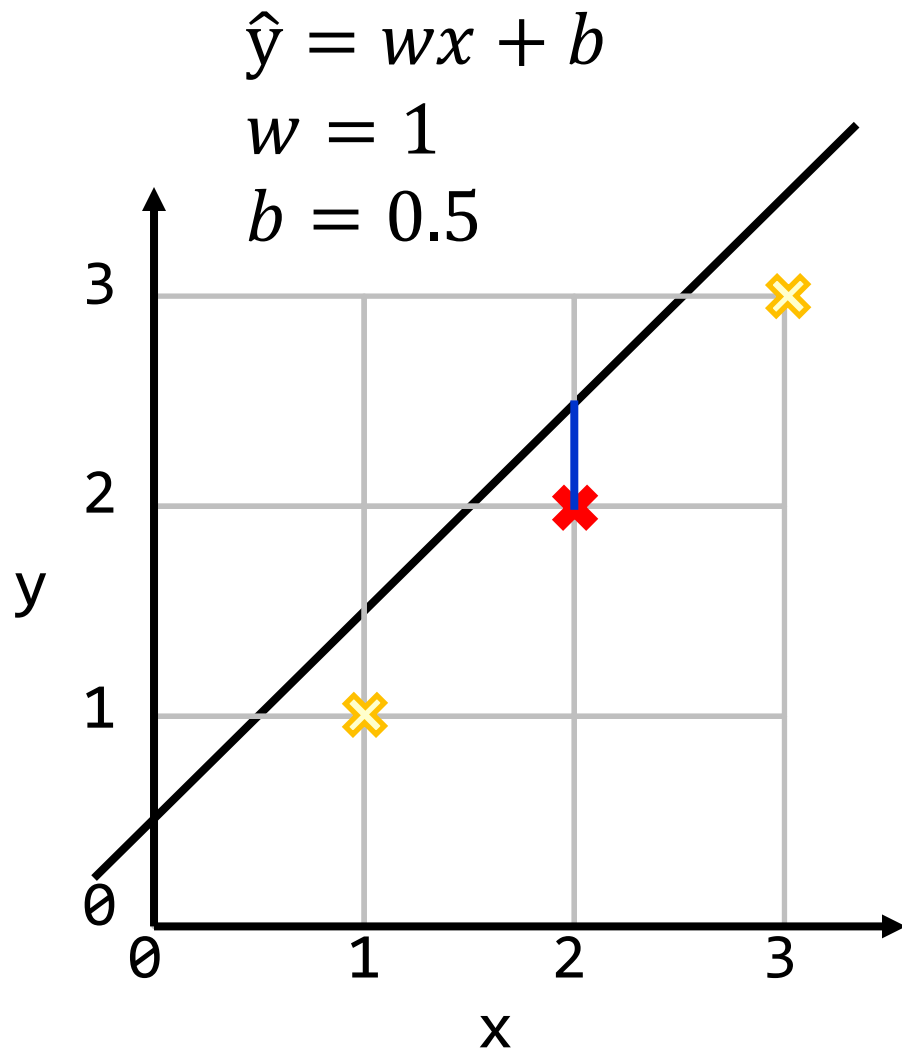
$$b = 0$$



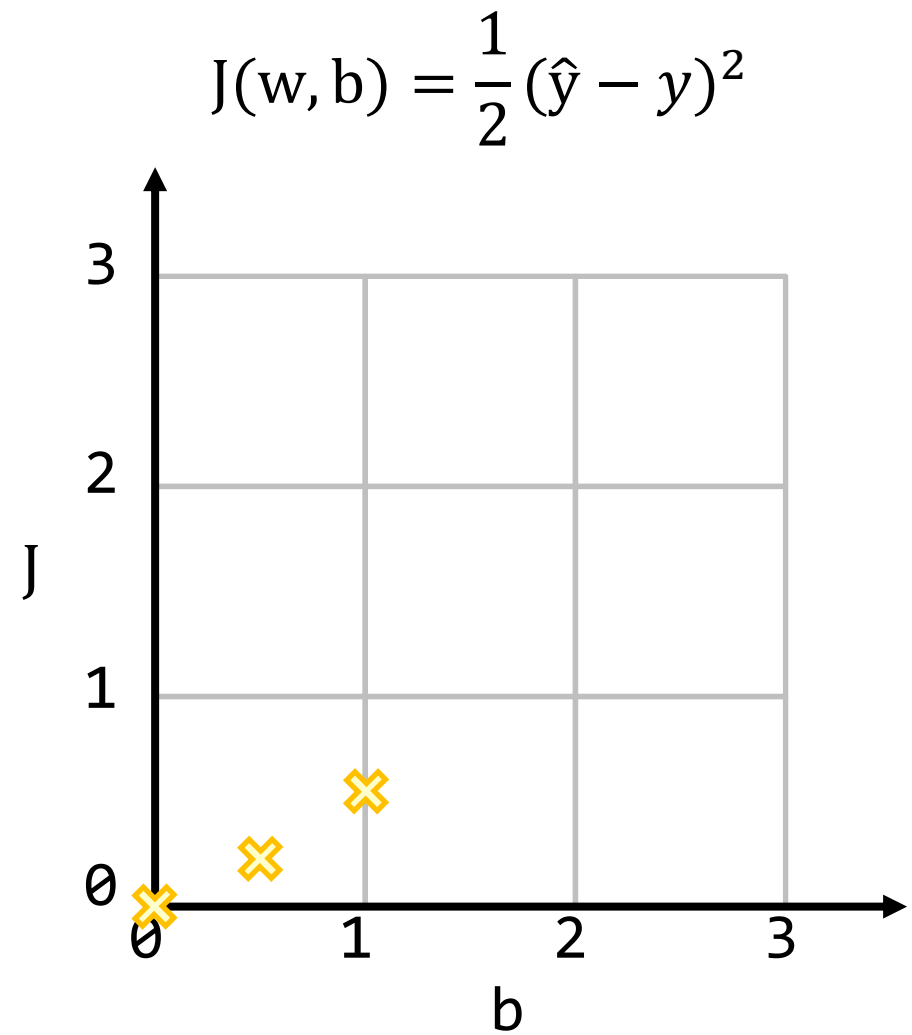
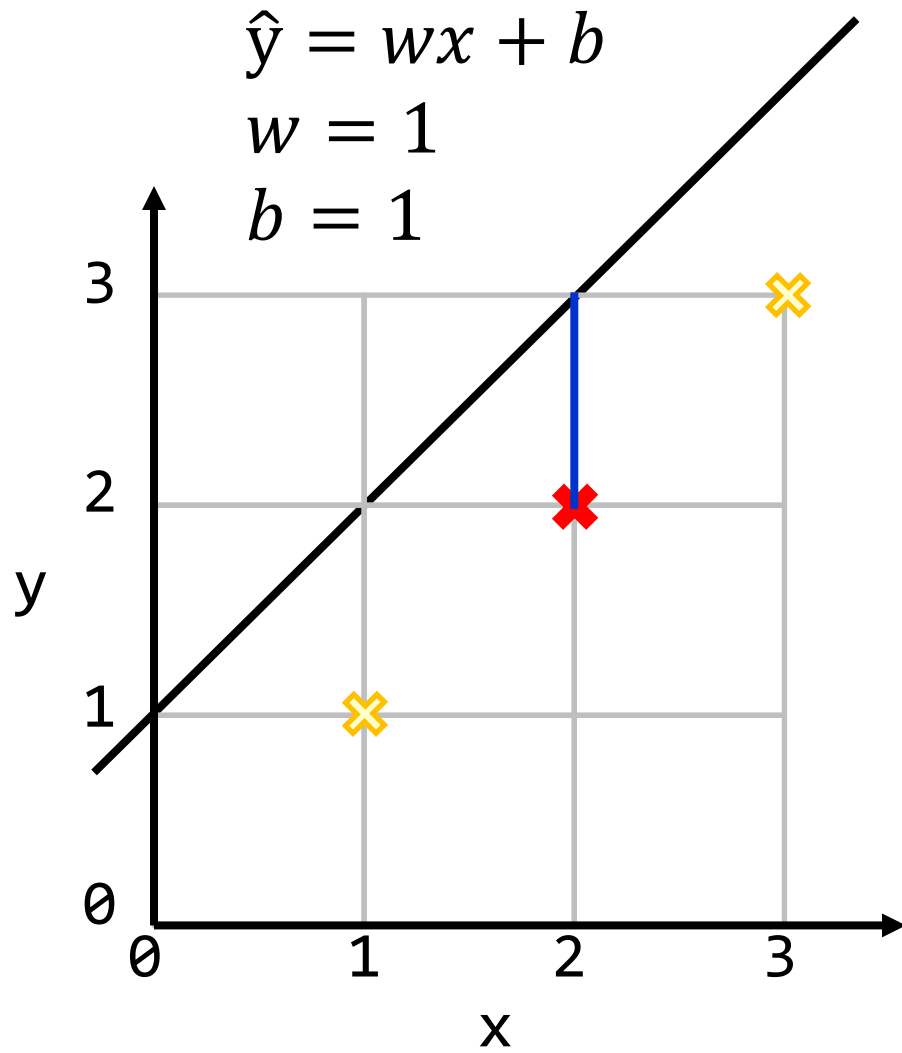
$$J(w, b) = \frac{1}{2} (\hat{y} - y)^2$$



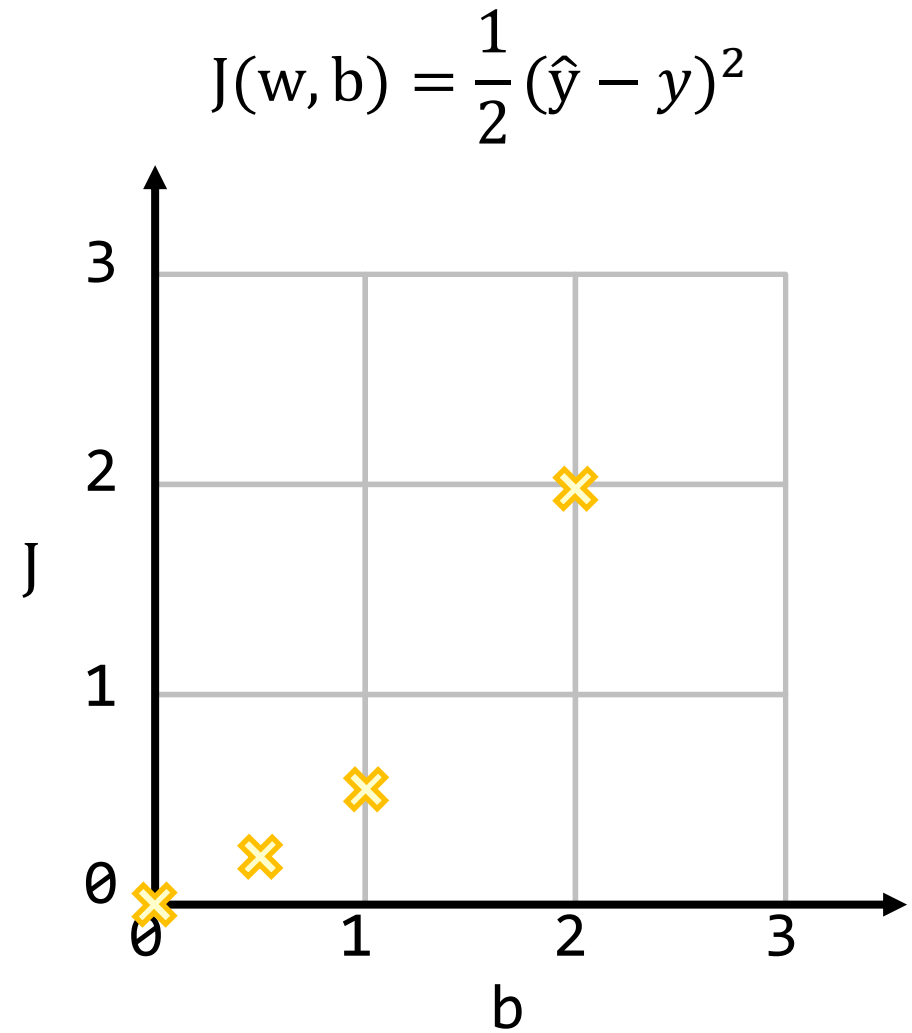
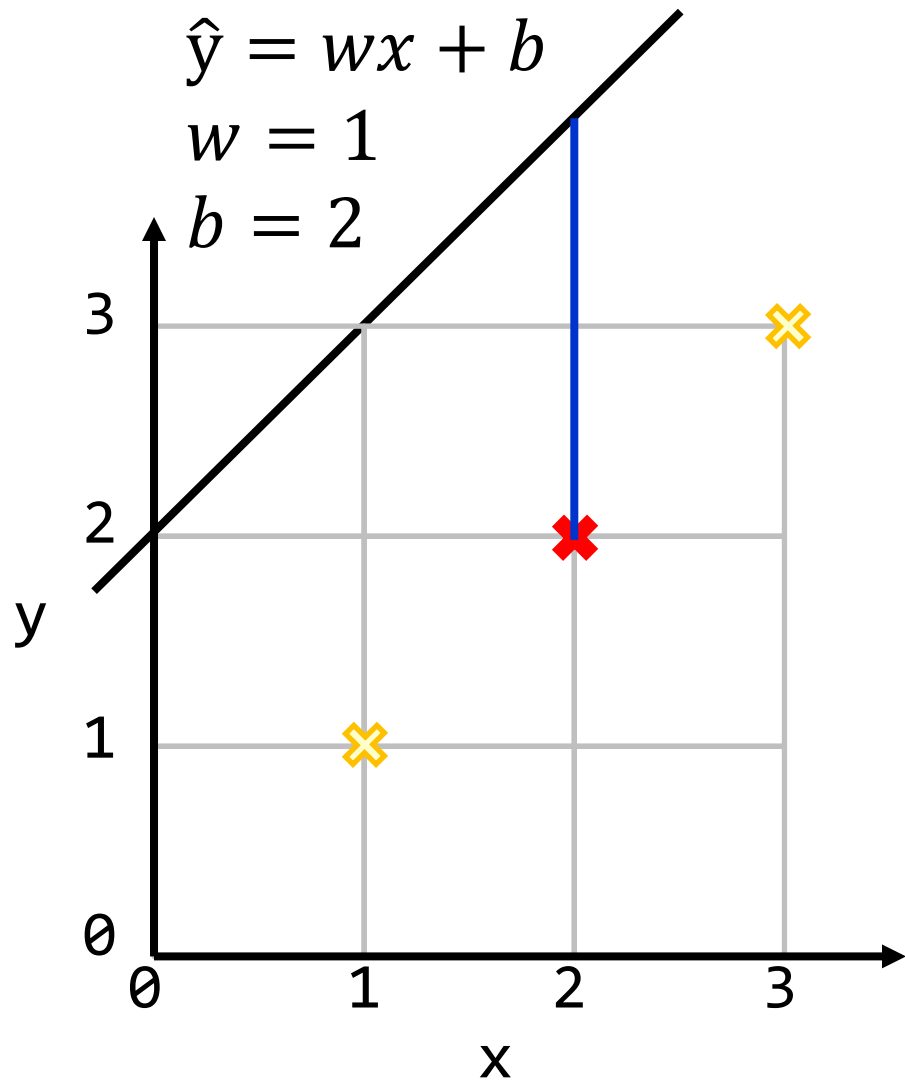
Cost function



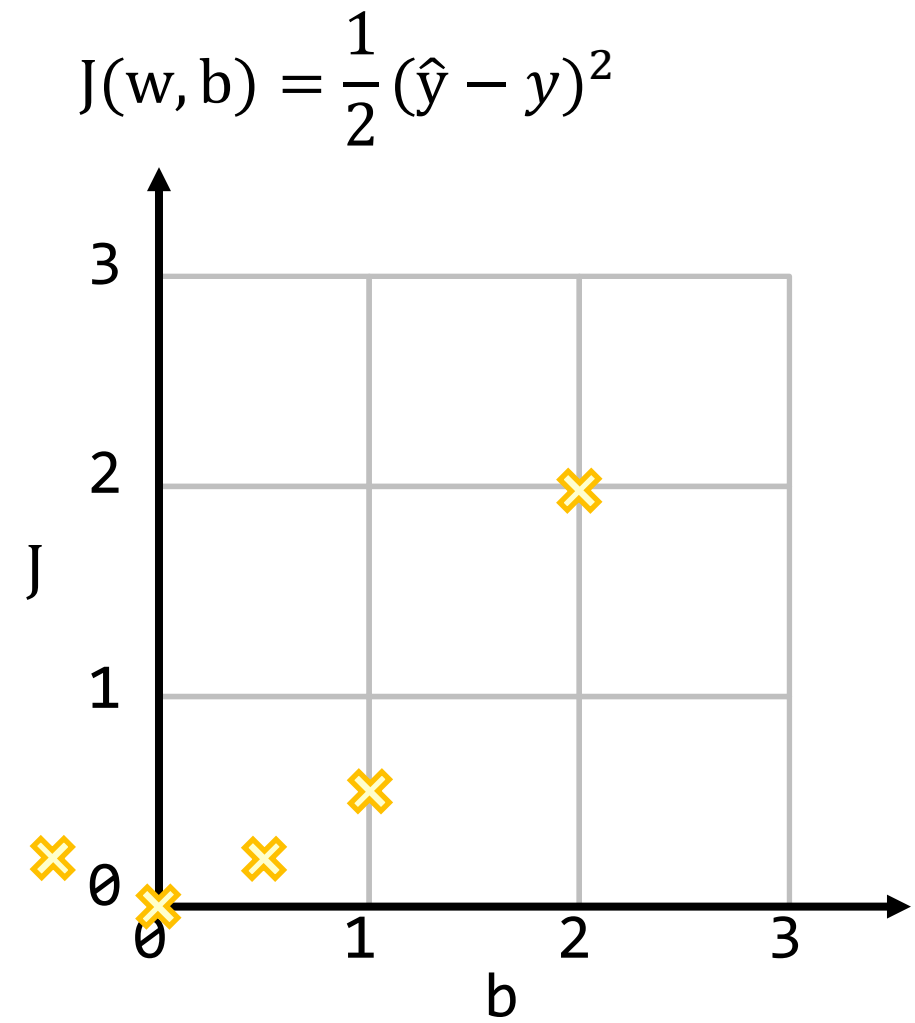
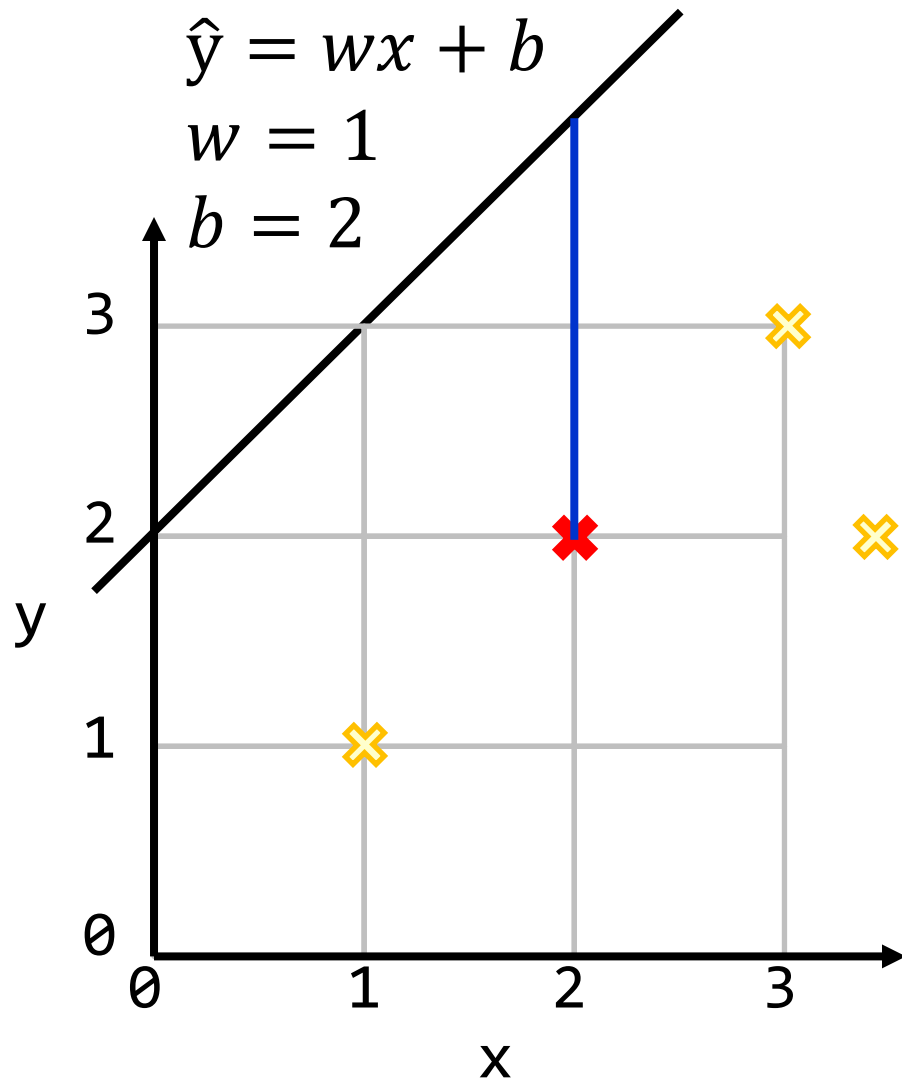
Cost function



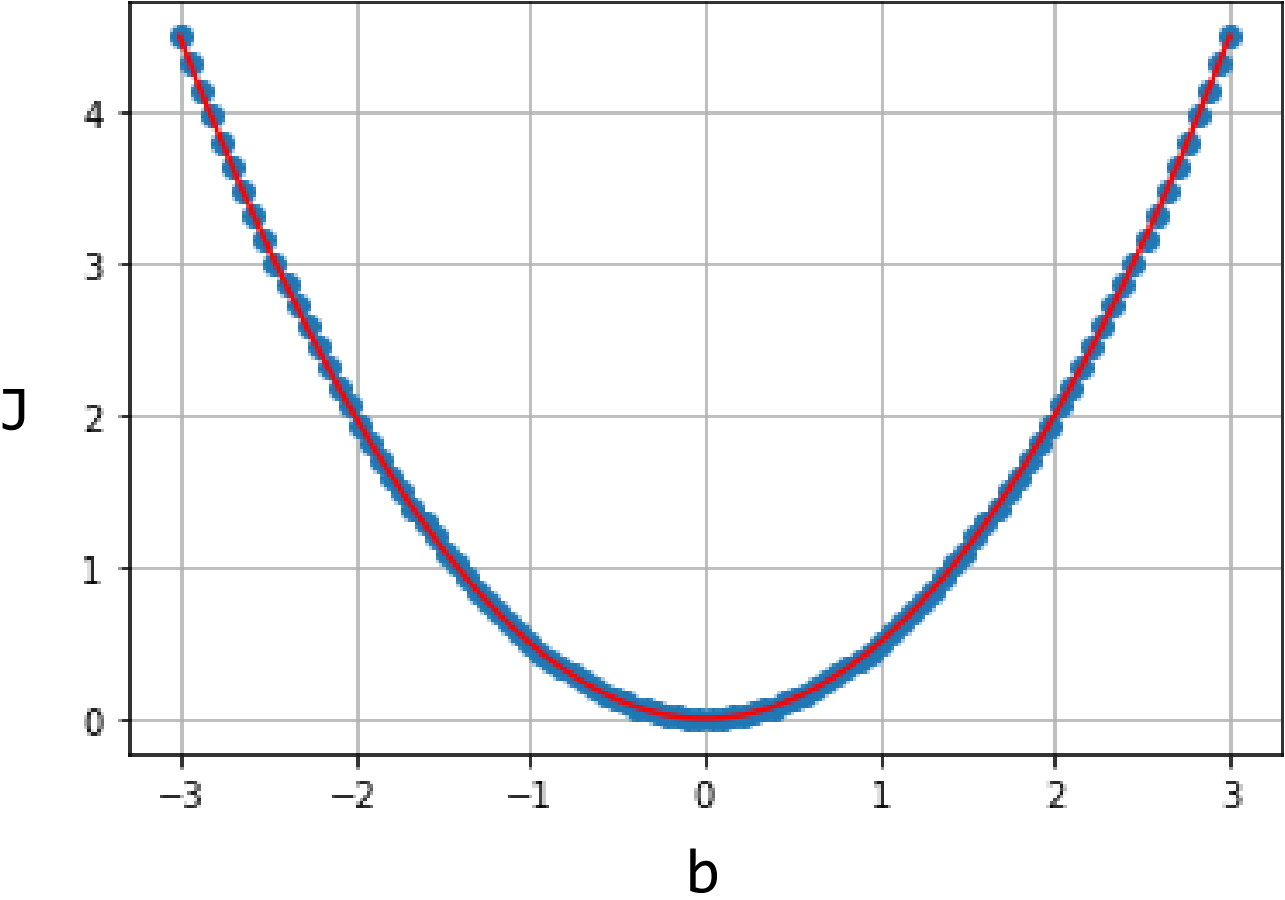
Cost function



Cost function



Cost function



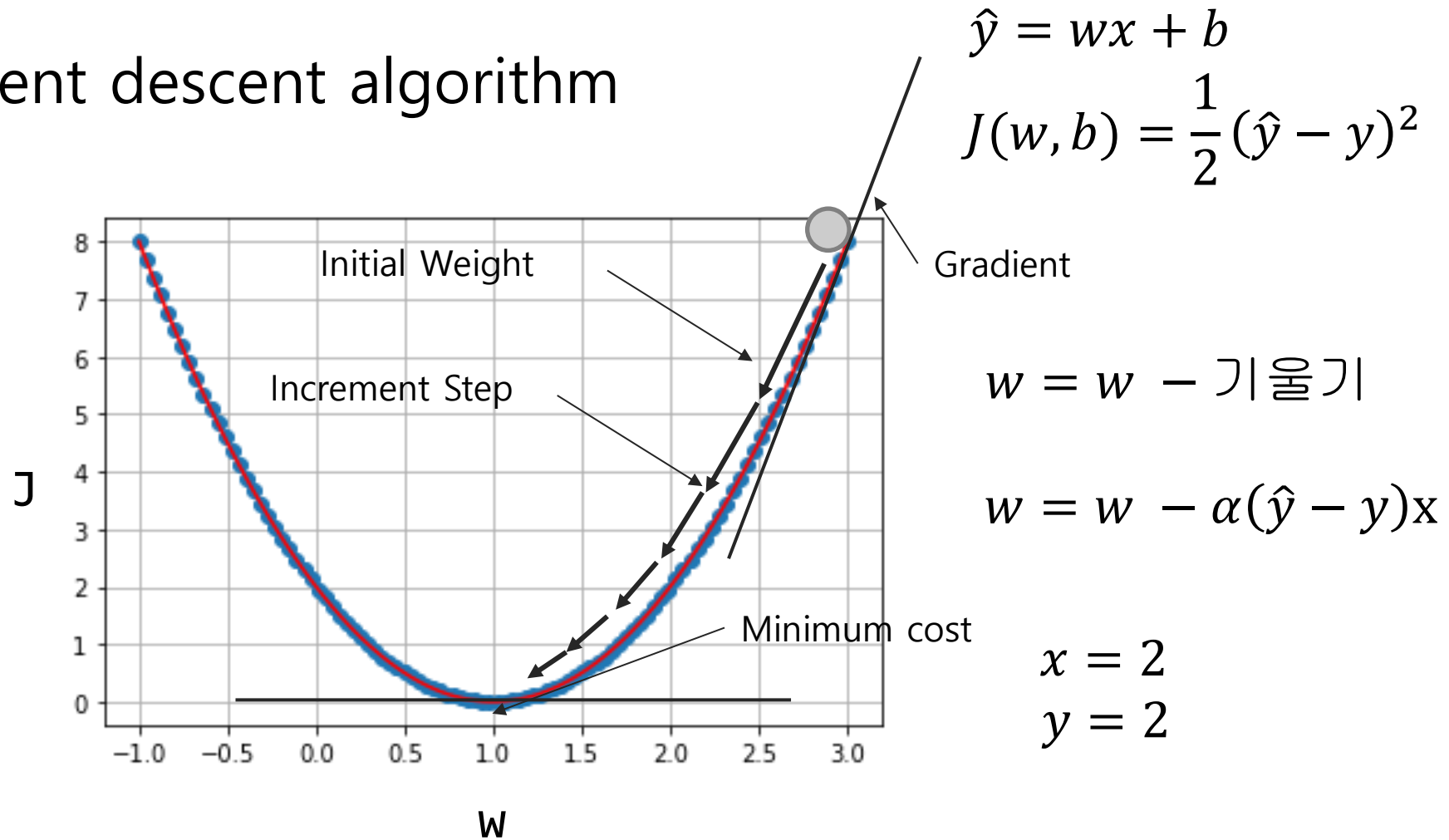
Gradient descent algorithm

- 비용 함수 최소화
- 경사 하강은 많은 최소화 문제에 사용된다.
- 주어진 비용 함수, 비용 (W, b) 에 대해 비용을 최소화하기 위해 W, b 를 찾는다.
- 일반적인 함수 : 비용 (w_1, w_2, \dots) 에 적용 가능

작동 방식

- 초기 추측으로 시작
 - 0,0 (또는 다른 값)에서 시작
 - W 와 b 를 약간 변경하여 $\text{cost}(W, b)$ 의 비용을 줄이려고 노력
- 매개 변수를 변경할 때마다 가능한 가장 낮은 $\text{cost}(W, b)$ 을 감소시키는 기울기를 선택
- 반복
- 최소한의 지역으로 수렴 할 때까지 수행

Gradient descent algorithm



Gradient descent algorithm

수렴까지 반복
{

$$w = w - \alpha \frac{\partial}{\partial w} J(w, b)$$

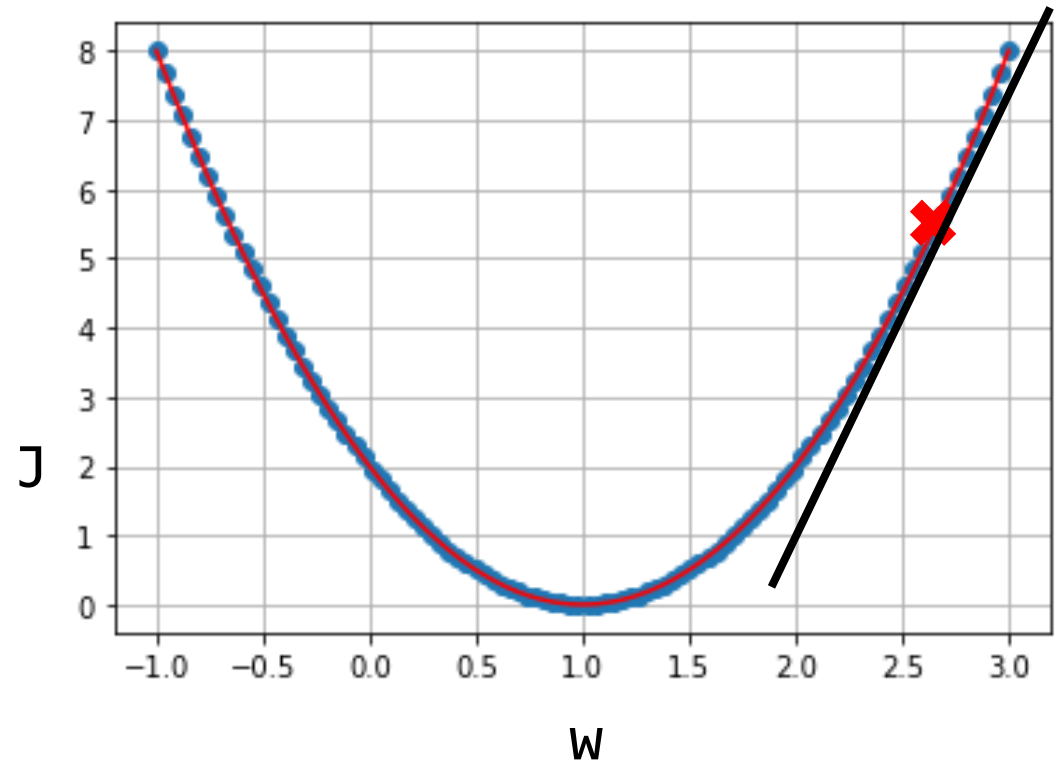
}

learning
rate

derivative

$$J(w, b) = \frac{1}{2} (\hat{y} - y)^2$$

$\frac{\partial}{\partial w} J(w, b)$ 는 $\frac{1}{2} (\hat{y} - y)^2$ 을 w 에 대해서 편미분 하면 된다.



Gradient descent algorithm

$$= \frac{\partial}{\partial w} \frac{1}{2} ((wx + b) - y)^2$$

$$\hat{y} = wx + b$$

$$= \frac{\partial}{\partial w} \frac{1}{2} ((wx + b)^2 - 2(wx + b)y + y^2)$$

$$J(w, b) = \frac{1}{2} (\hat{y} - y)^2$$

$$= \frac{\partial}{\partial w} \frac{1}{2} ((wx + b)^2 - 2ywx - 2by + y^2)$$

$$= \frac{\partial}{\partial w} \frac{1}{2} (w^2x^2 + 2wxb + \cancel{b^2} - 2ywx - \cancel{2by} + \cancel{y^2})$$

$$= \frac{1}{2} (2wx^2 + 2xb - 2yx)$$

$$= x(wx + b - y)$$

$$= x(\hat{y} - y)$$

$$\frac{\partial}{\partial \hat{y}} \frac{1}{2} (\hat{y} - y)^2 = \hat{y} - y$$

$$\frac{\partial}{\partial w} (wx + b) = x$$

$$\frac{\partial J}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial w} = x(\hat{y} - y)$$

Chain Rule 사용

$$\hat{y} = wx + b$$

$$J(w, b) = \frac{1}{2} (\hat{y} - y)^2$$

$$\frac{\partial}{\partial w} J(w, b) = \frac{\partial J}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial w}$$

Gradient descent algorithm

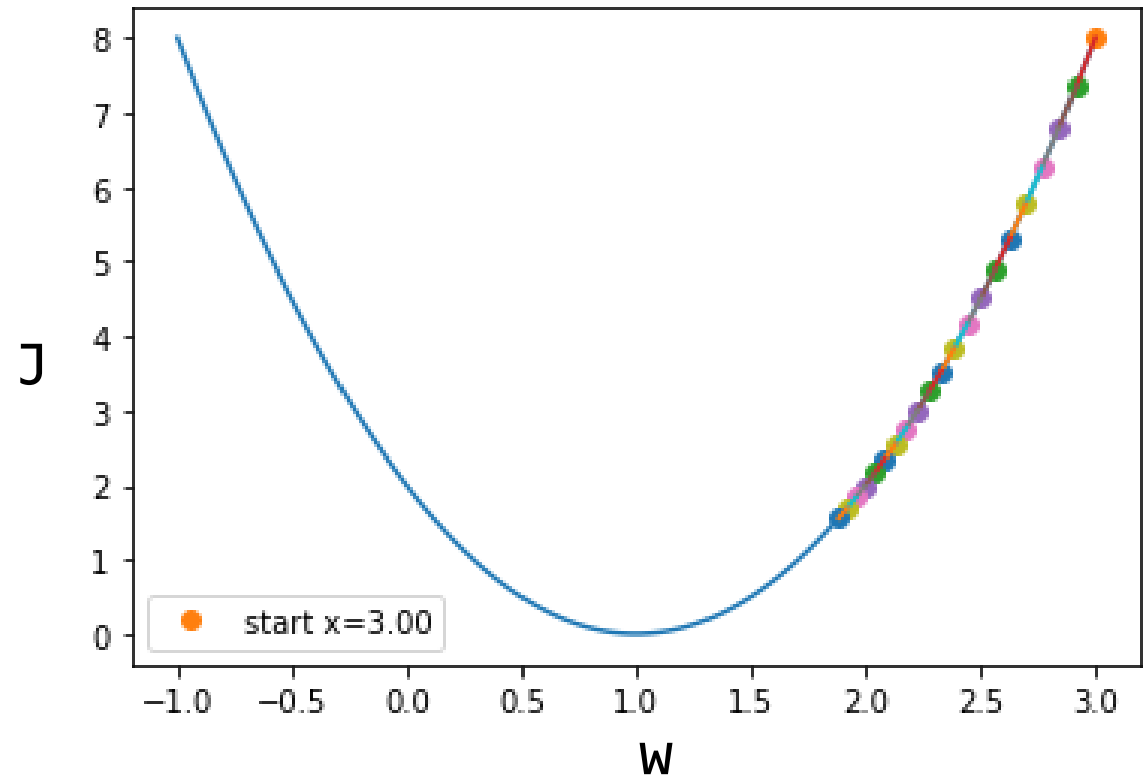
수렴까지 반복

{

$$w = w - \alpha * (\hat{y} - y) * x$$

}

α 가 0.01인 경우



Gradient descent algorithm

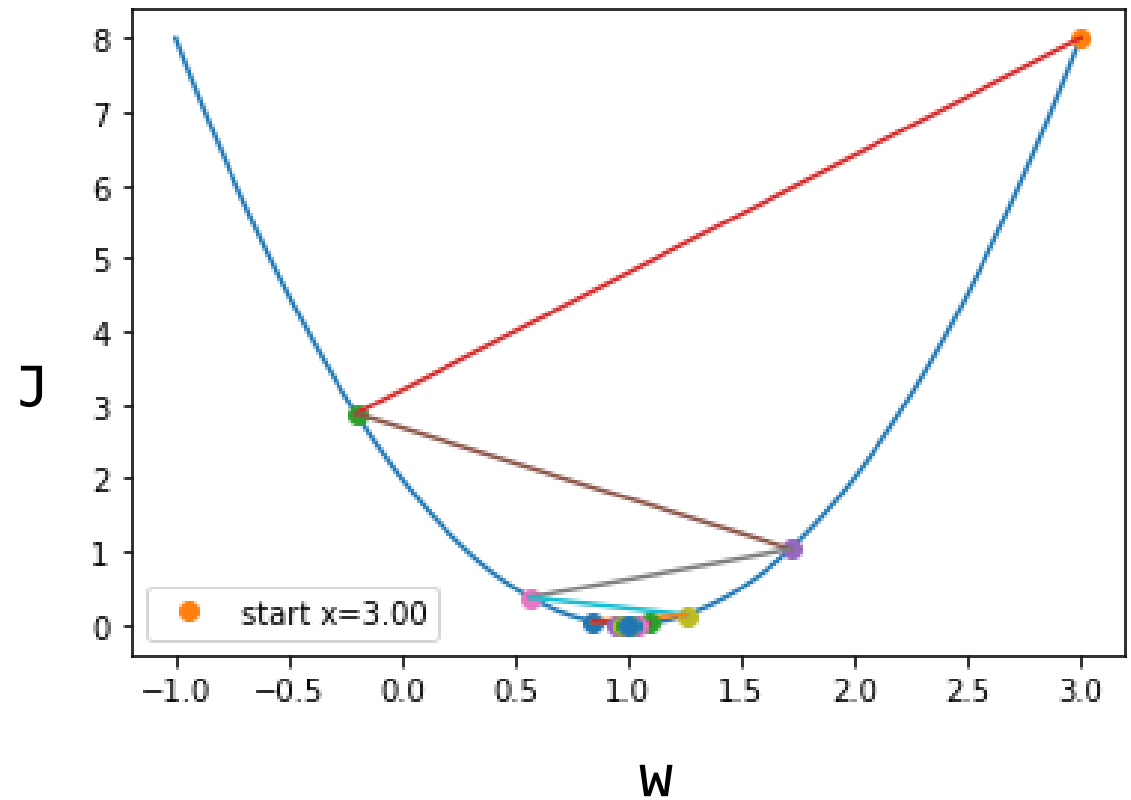
수렴까지 반복

{

$$w = w - \alpha * (\hat{y} - y) * x$$

}

α 가 **0.4**인 경우



Gradient descent algorithm

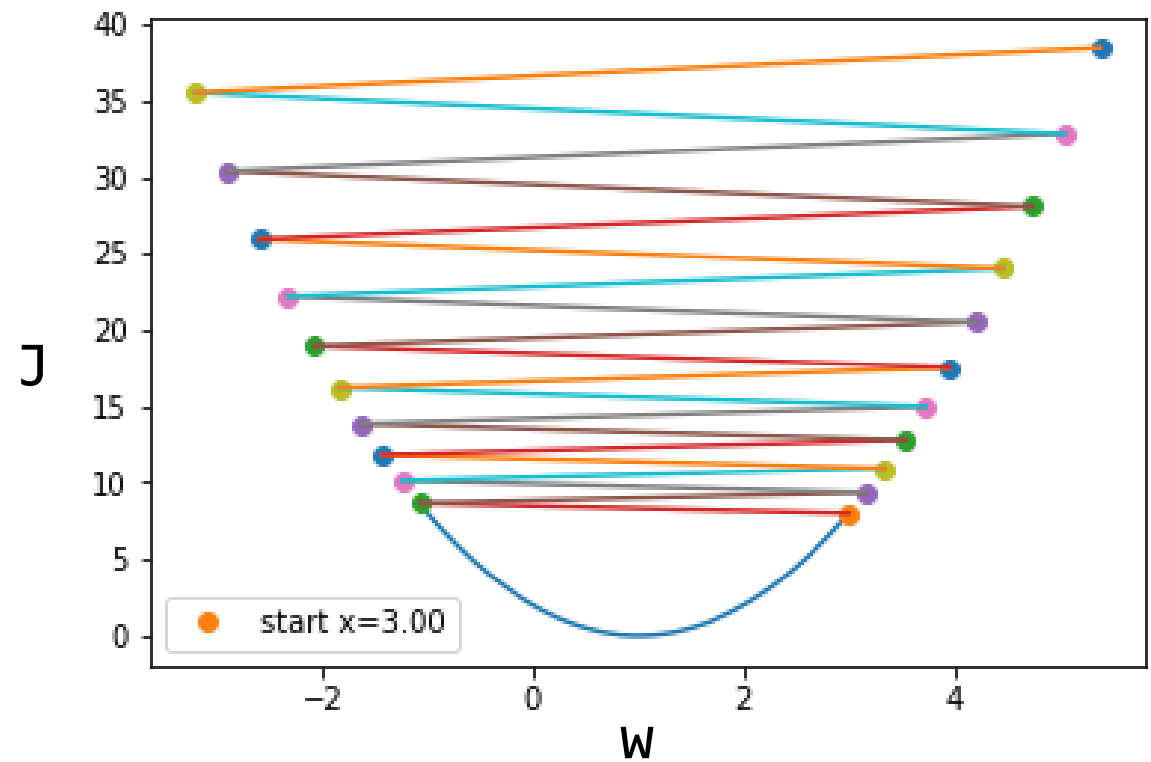
수렴까지 반복

{

$$w = w - \alpha * (\hat{y} - y) * x$$

}

α 가 0.51인 경우



$$\frac{\partial J}{\partial \hat{y}} = \frac{1}{2} (\hat{y} - y)^2 = \hat{y} - y$$

$$\frac{\partial \hat{y}}{\partial w} = wx + b = 1$$

$$\frac{\partial J}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial w} = (\hat{y} - y)$$

Chain Rule 사용

$$\hat{y} = wx + b$$

$$J(w,b) = \frac{1}{2} (\hat{y} - y)^2$$

$$\frac{\partial}{\partial b} J(w,b) = \frac{\partial J}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial b}$$

Gradient descent algorithm

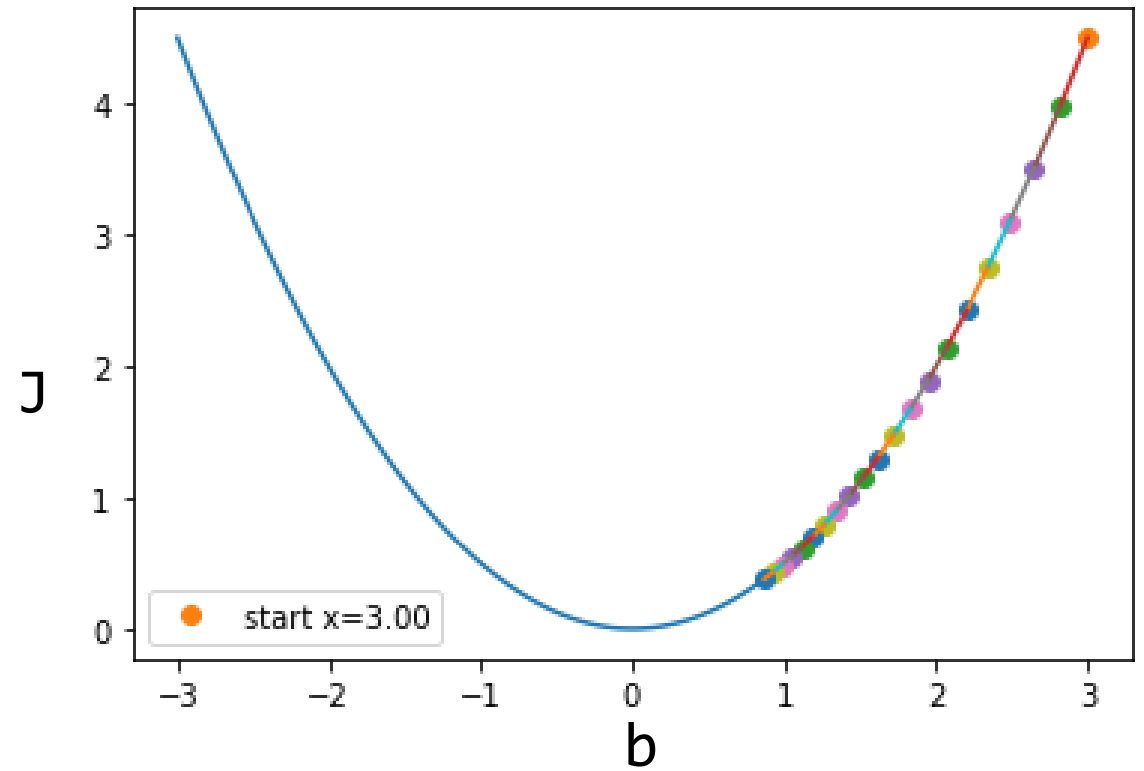
수렴까지 반복

{

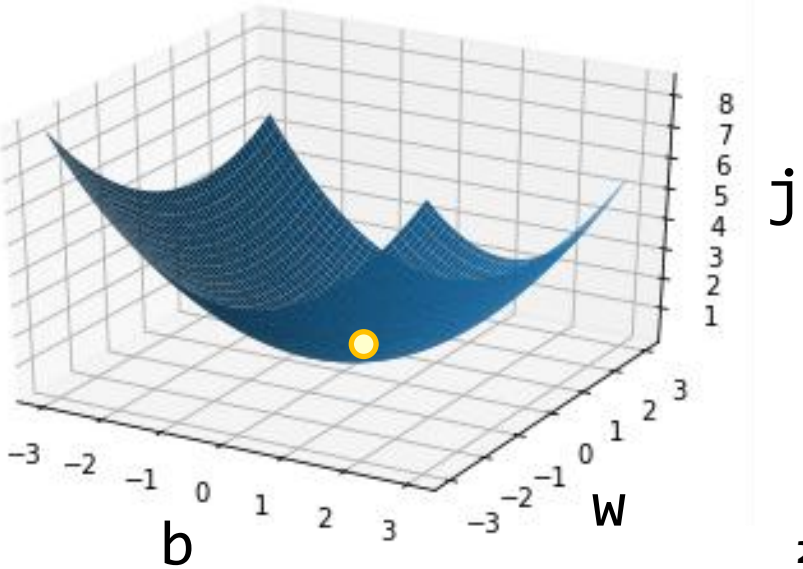
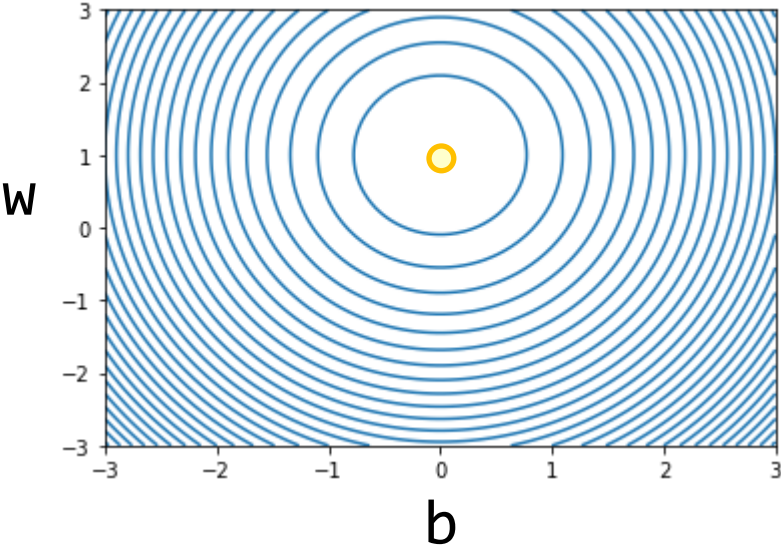
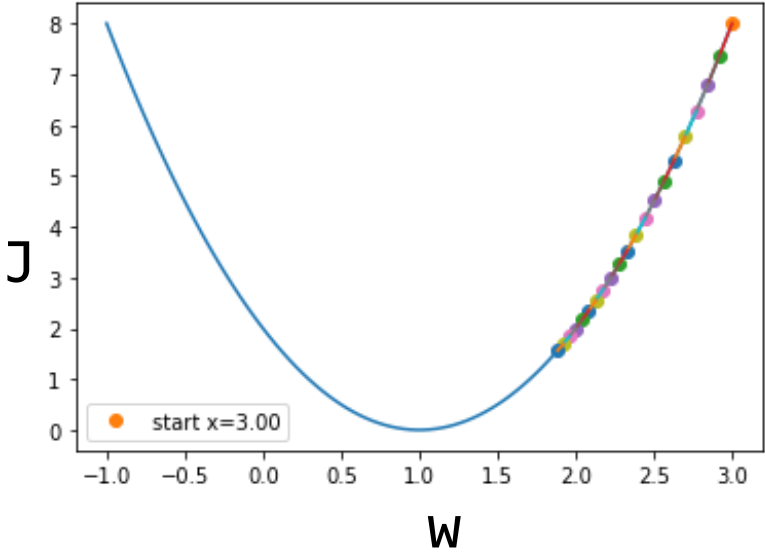
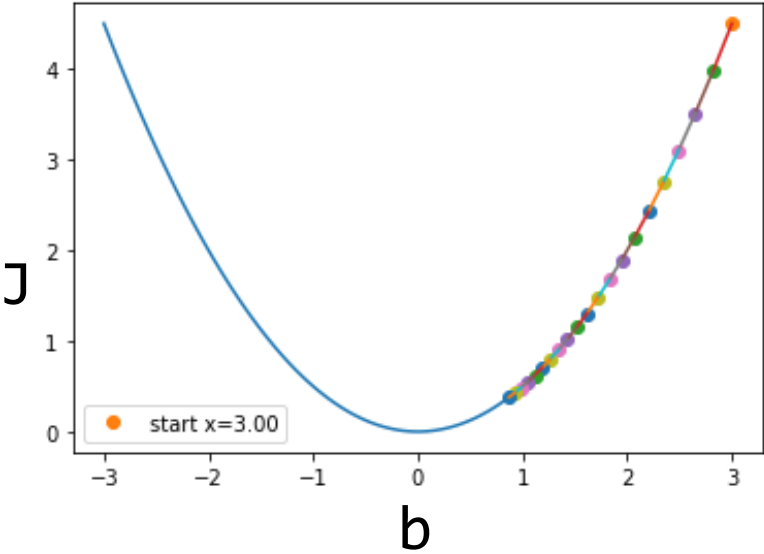
$$b = b - \alpha * (\hat{y} - y)$$

}

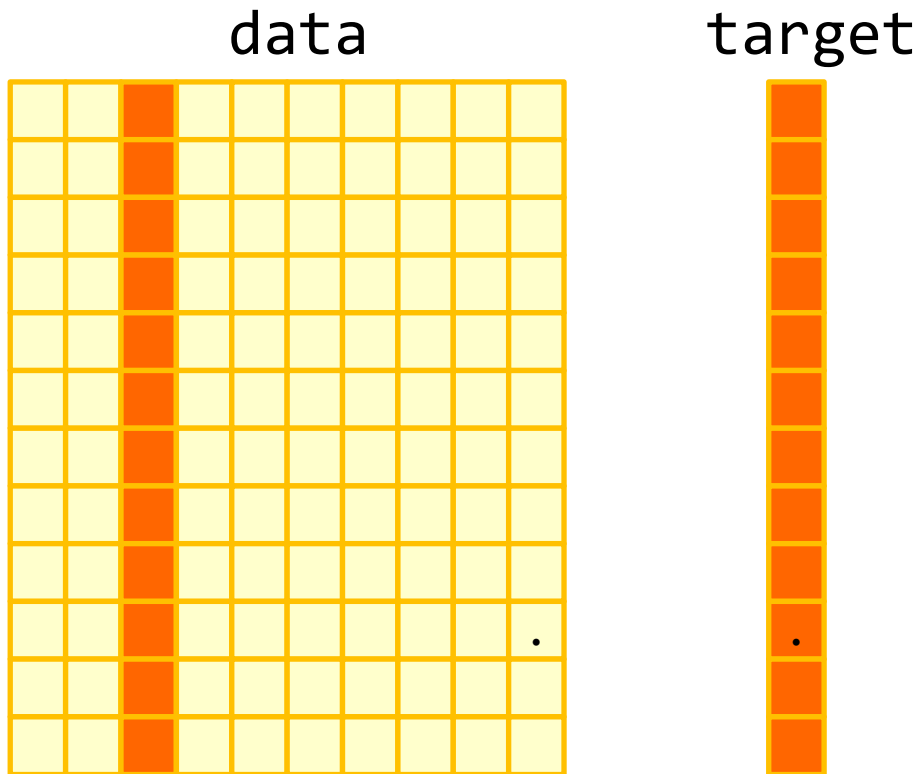
α 가 **0.03**인 경우



Gradient descent algorithm



data set 준비(sklearn 이용)

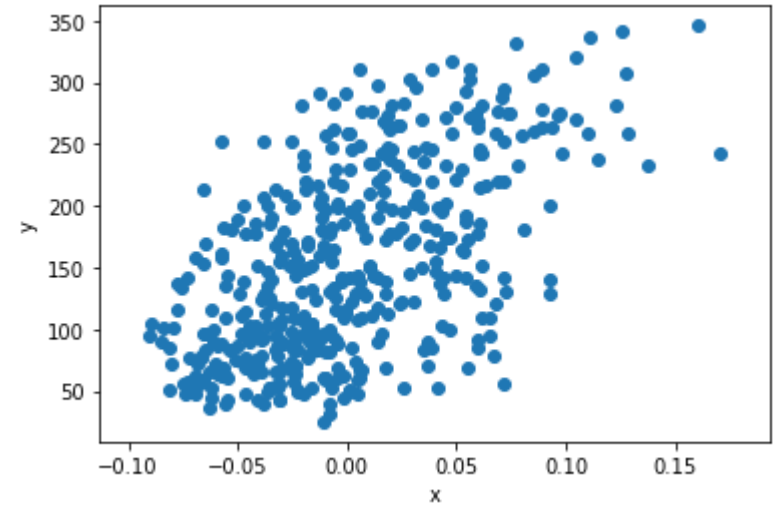


data set 준비(sklearn 이용)

```
from sklearn.datasets import load_diabetes
diabetes = load_diabetes()
import matplotlib.pyplot as plt

x = diabetes.data[:, 2]
y = diabetes.target

plt.scatter(x, y)
plt.xlabel('x')
plt.ylabel('y')
plt.show()
```



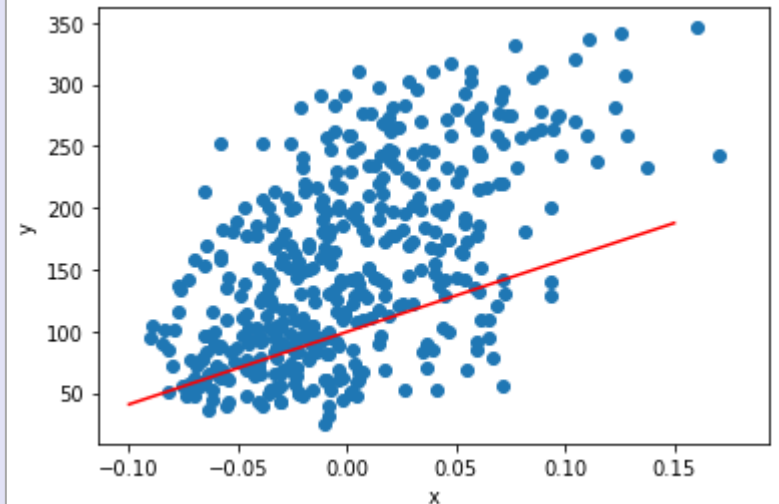
역전파 구현 I

역전파 구현

1. batch=1
2. 점의 전체 개수만큼 1번 순회 (epoch=1)
: SGD=>원소 하나마다 w, b 갱신

```
w = 1.0
b = 1.0
rate = 0.01
for x_i, y_i in zip(x, y):
    y_hat = x_i * w + b
    err = y_hat - y_i
    w = w - rate * err * x_i
    b = b - rate * err

plt.scatter(x, y)
pt1 = (-0.1, -0.1 * w + b)
pt2 = (0.15, 0.15 * w + b)
plt.plot([pt1[0], pt2[0]], [pt1[1], pt2[1]], 'r-')
plt.xlabel('x')
plt.ylabel('y')
plt.show()
```

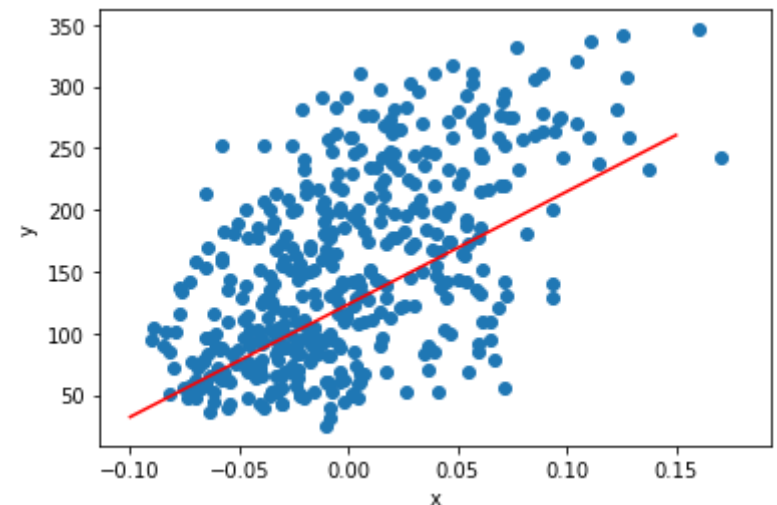


역전파 구현 II

역전파 구현

1. batch=1
2. epoch=100번 만큼 외부 루프에서 순회
3. 점의 전체 개수만큼 1번 순회
: SGD=>원소 하나마다 w, b 갱신

```
for i in range(1, 100):  
    for x_i, y_i in zip(x, y):  
        y_hat = x_i * w + b  
        err = y_i - y_hat  
        w_rate = x_i  
        w = w + w_rate * err  
        b = b + 1 * err  
plt.scatter(x, y)  
pt1 = (-0.1, -0.1 * w + b)  
pt2 = (0.15, 0.15 * w + b)  
plt.plot([pt1[0], pt2[0]], [pt1[1], pt2[1]], 'r-')  
plt.xlabel('x')  
plt.ylabel('y')  
plt.show()
```



Logistic Regression

- 로지스틱 회귀 란?
 - Classification(분류)
 - Logistic vs Linear
- 동작 방식
 - 가설 표현
 - Sigmoid 함수
 - Decision Boundary(결정경계)
 - Cost Function
 - Optimizer (Gradient Descent)

Classification

Binary Classification(이진 분류) 란?

: 값은 0 또는 1

- Exam : Pass or **Fail**
- Spam : Not Spam or **Spam**
- Face : Real or **Fake**
- Tumor : Not Malignant or **Malignant**

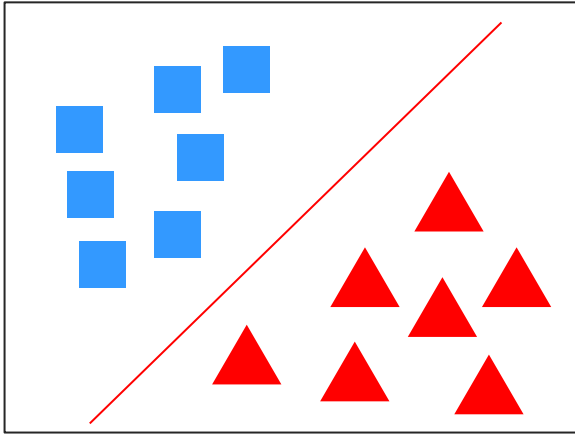
머신 러닝으로 시작하려면 변수를 인코딩해야 한다.

```
x_train = [[1,2], [2,3], [3,1], [4,3], [5,3], [6,2]]
```

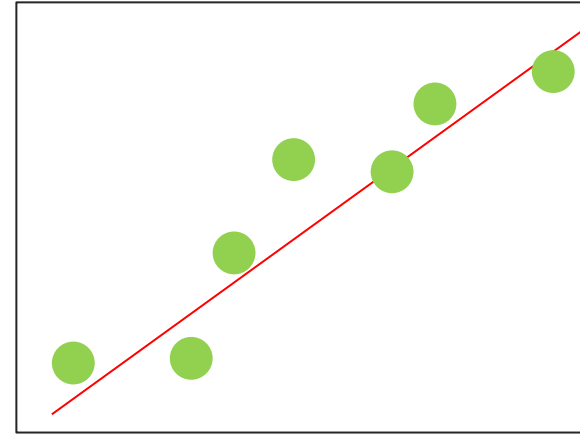
```
y_train = [[0], [0], [0], [1], [1], [1]]
```

Logistic vs Linear

로지스틱 회귀와 선형 회귀의 차이점은?

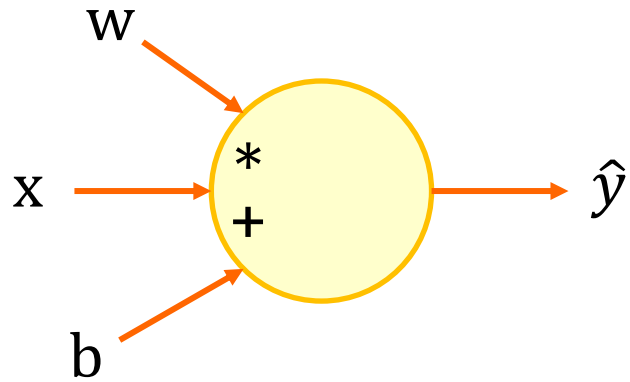


Discrete(분리) : 분류 목적
신발 사이즈 / 회사의 근로자수

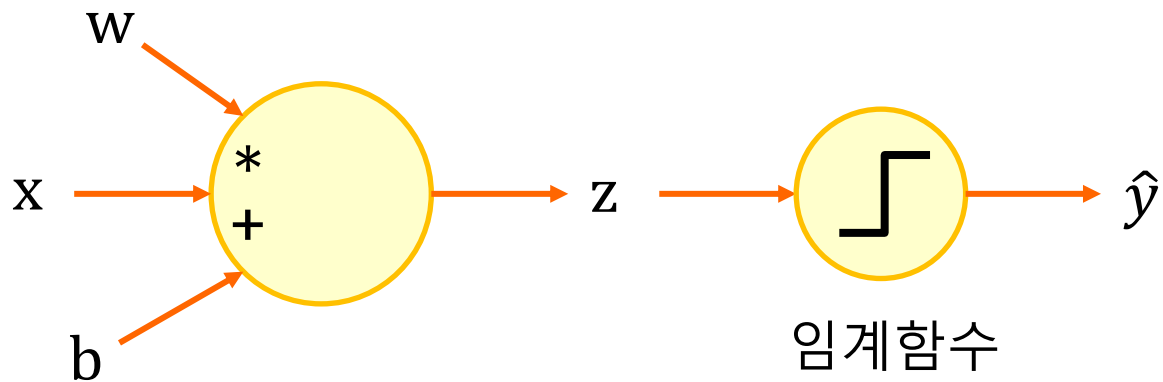


Continuous (지속적인) : 값의 예측
시간/ 무게 / 높이

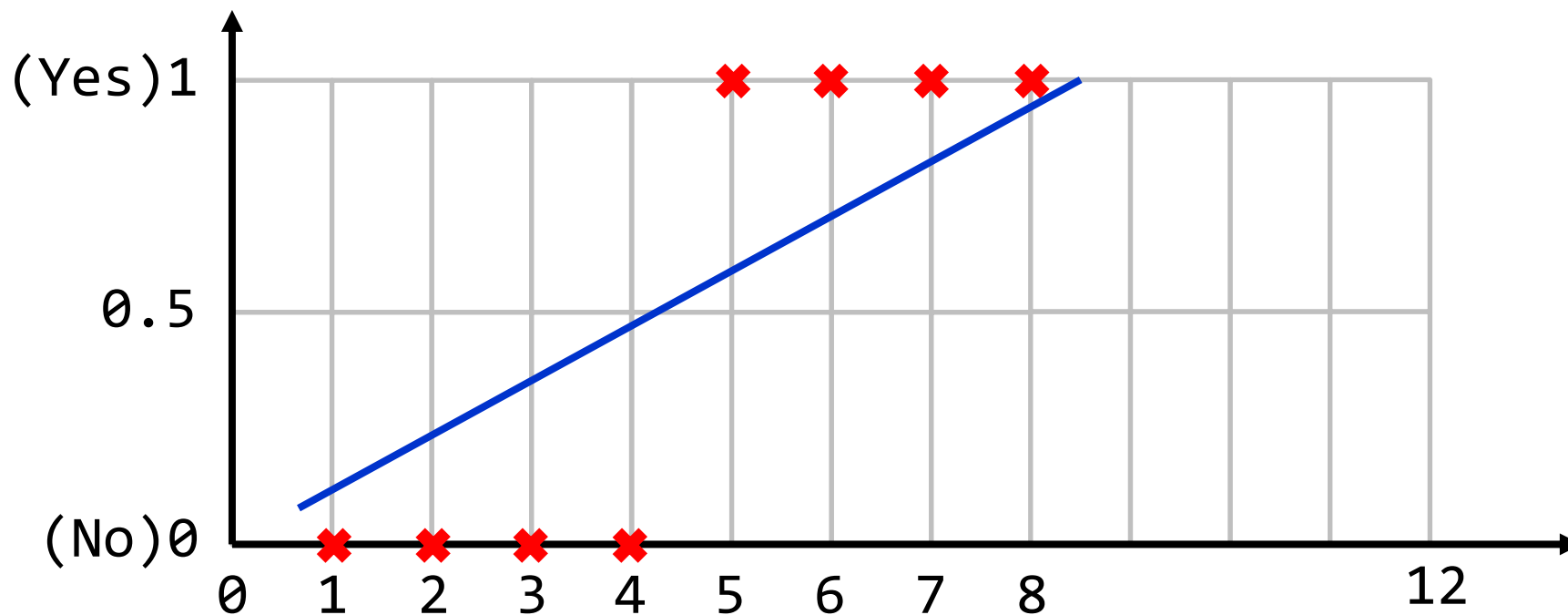
퍼셉트론



퍼셉트론



Sigmoid 함수 사용 이유



Threshold classifier output at 0.5:

If $\hat{y} \geq 0.5$, predict "y=1"

If $\hat{y} < 0.5$, predict "y=0"

1 => 양성

2 => 양성

3 => 양성

4 => 양성

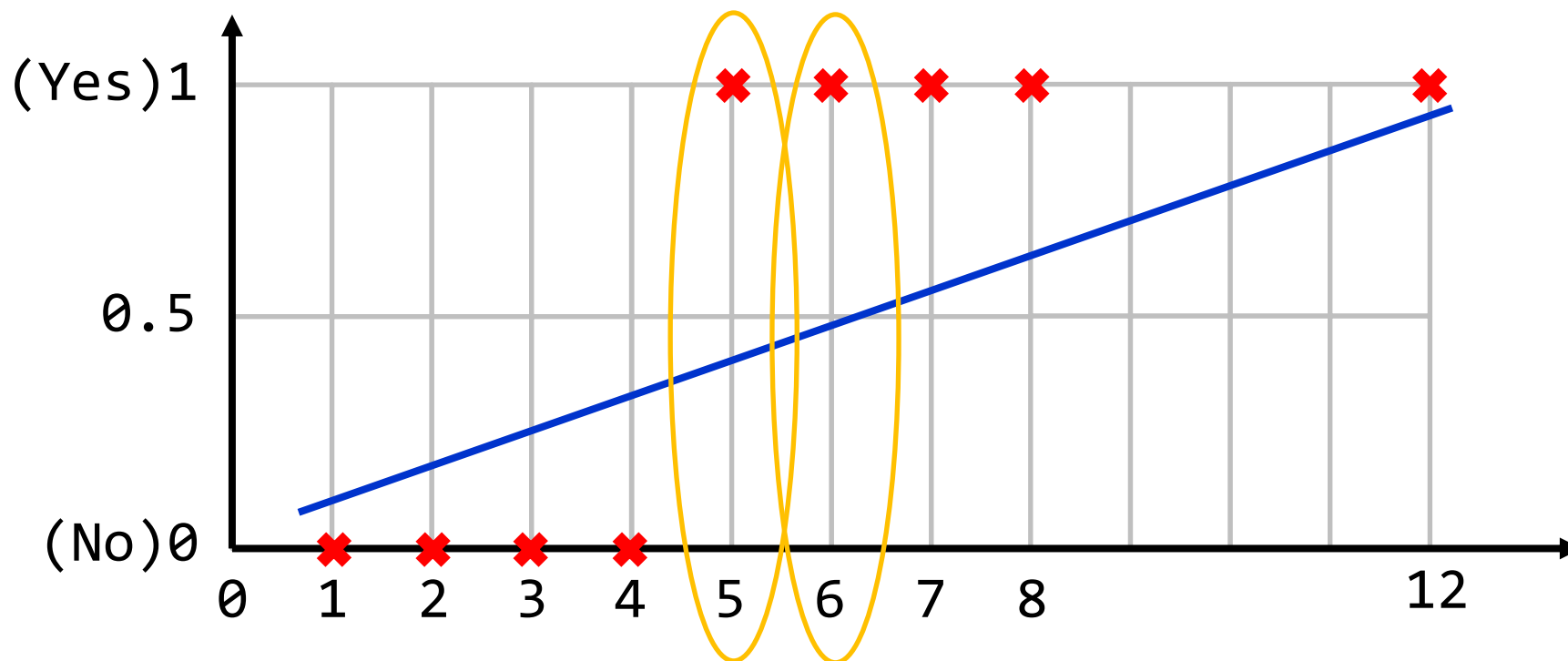
5 => 악성

6 => 악성

7 => 악성

8 => 악성

Sigmoid 함수 사용 이유



Threshold classifier output at 0.5:

If $\hat{y} \geq 0.5$, predict "y=1"

If $\hat{y} < 0.5$, predict "y=0"

1 => 양성

2 => 양성

3 => 양성

4 => 양성

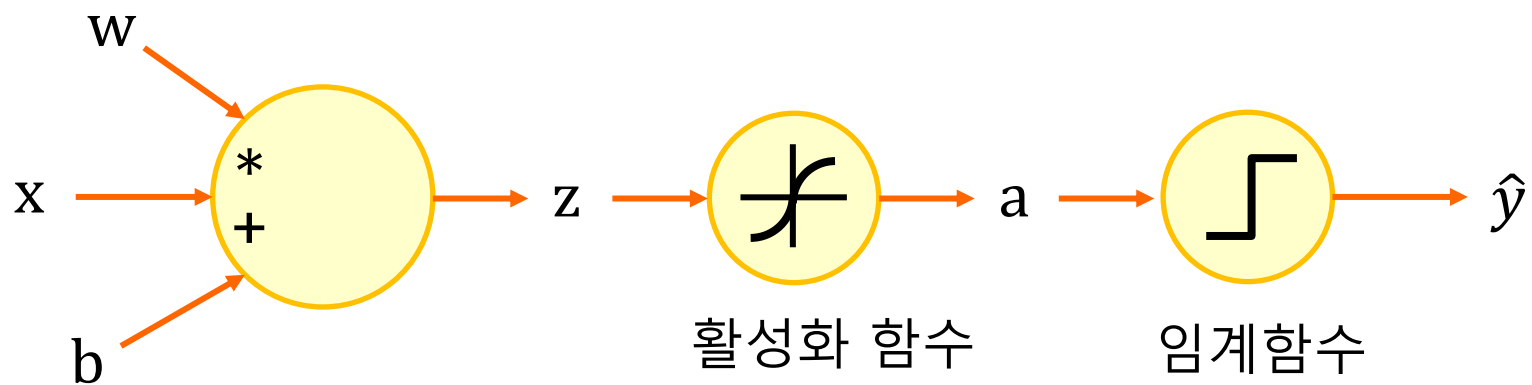
5 => 양성

6 => 양성

7 => 악성

8 => 악성

Sigmoid 함수 사용 이유



Classification: $y=0$ or 1

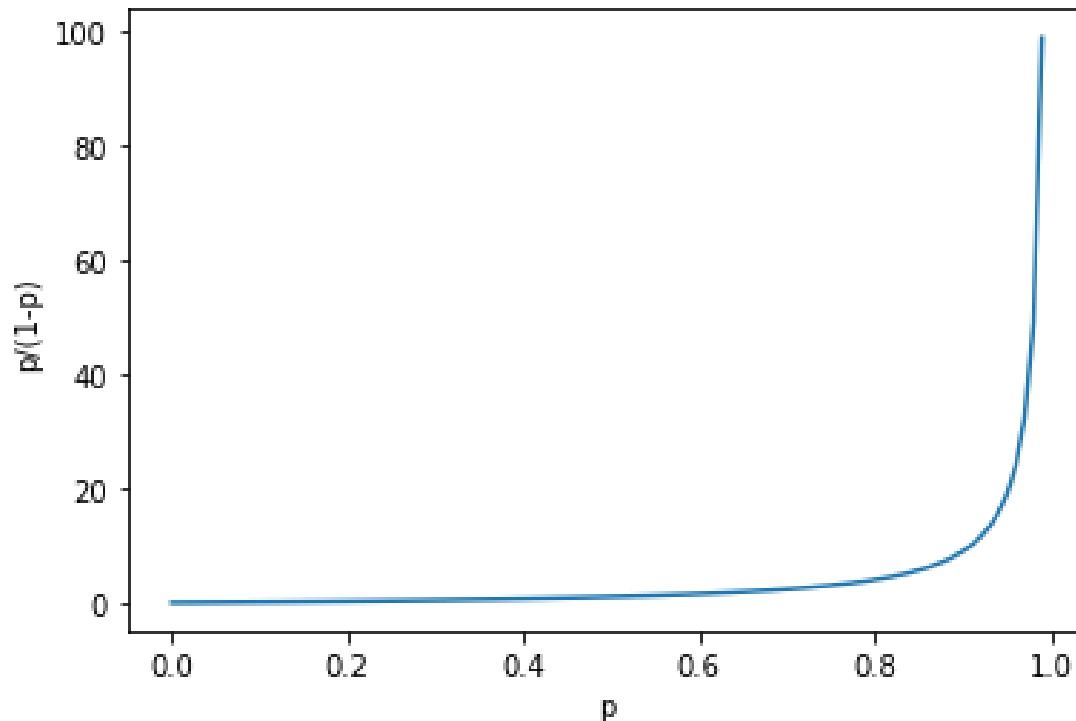
\hat{y} can be > 1 or < 0

Logistic Regression: $0 < \hat{y} < 1$

분류 문제 사용

$$\text{OR(odds ratio)} = \frac{p}{1-p} \quad (p=\text{성공확률})$$

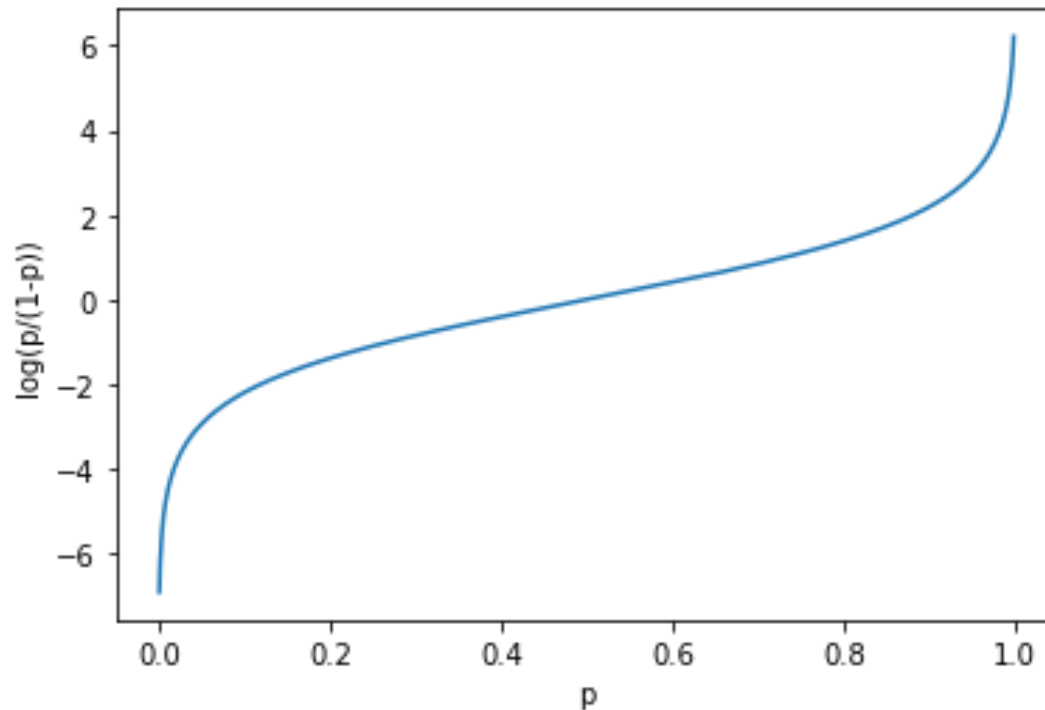
오즈 비를 그래프로 그리면 다음과 같다.
p가 0부터 1까지 증가할 때 오즈 비의 값은 처음에는
천천히 증가하지만 p가 1에 가까워지면 급격히 증가한다.



로짓 함수(Logit function)

$$\text{logit}(p) = \log\left(\frac{p}{1-p}\right) \quad (p=\text{성공 확률})$$

로짓 함수는 p 가 0.5일 때 0이 되고 p 가 0과 1일 때 각각 무한대로 음수와 양수가 되는 특징을 가진다.



시그모이드 함수(Sigmoid function)

$$\text{logit}(p) = \log\left(\frac{p}{1-p}\right) = z$$

로지틱 함수의 유도 : p에 대해 정리

$$\log\left(\frac{p}{1-p}\right) = z$$

$$e^{\log\left(\frac{p}{1-p}\right)} = e^z$$

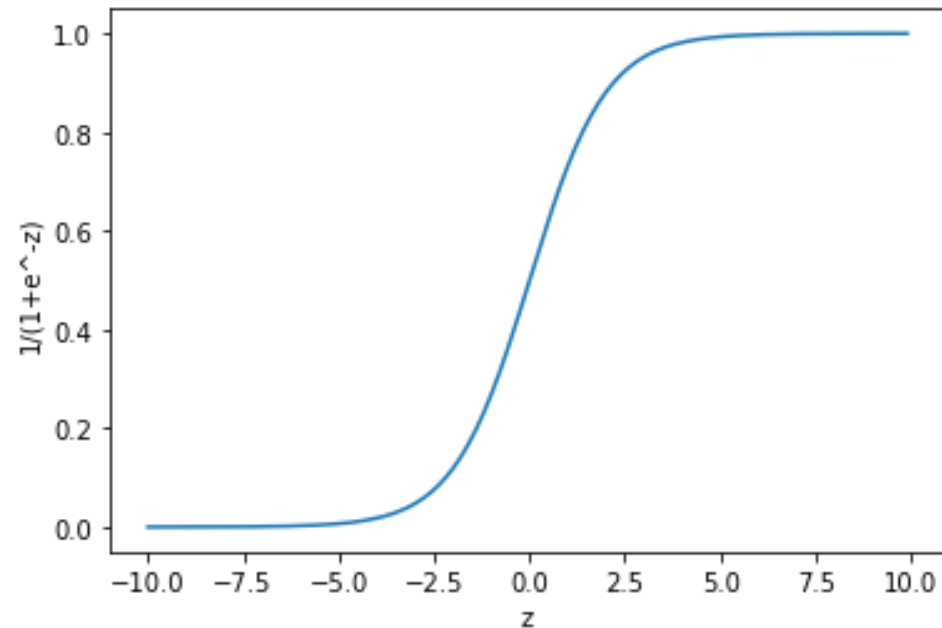
$$\frac{p}{1-p} = e^z$$

$$p = (1 - p) * e^z$$

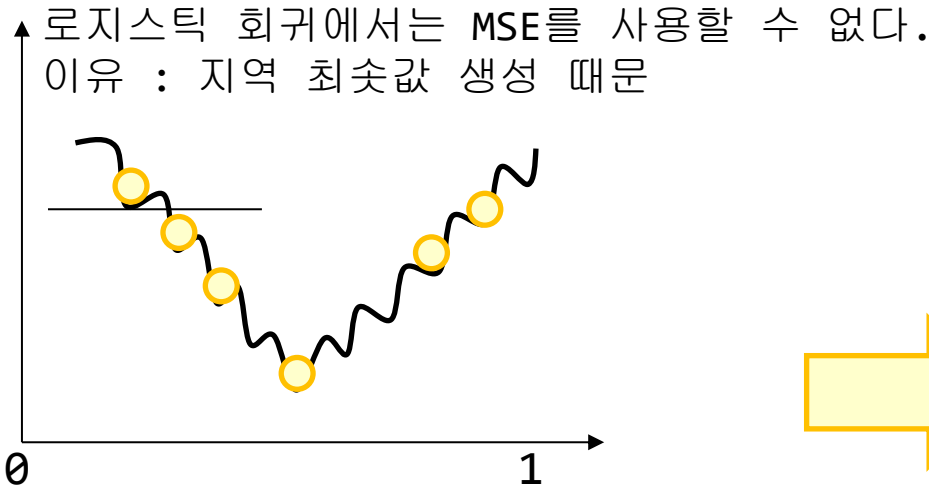
$$p = e^z - p * e^z$$

$$p + p * e^z = e^z$$

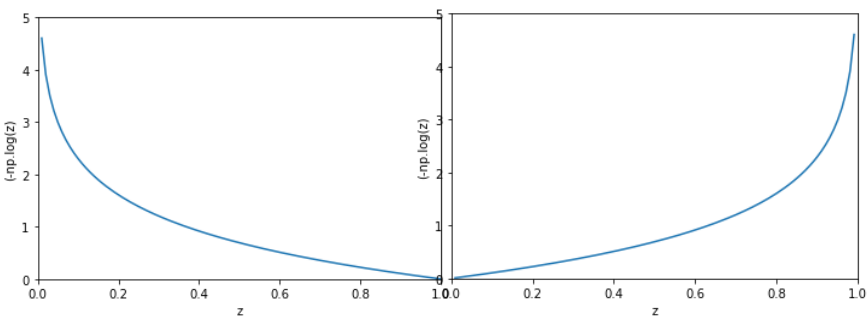
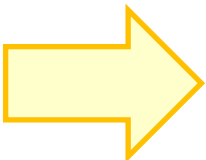
$$p = \frac{e^z}{1+e^z} = \frac{1}{1+e^{-z}}$$



A convex logistic regression cost function



Binary Cross Entropy 함수



$$\text{sigmoid} - \text{actual} = \text{error}$$

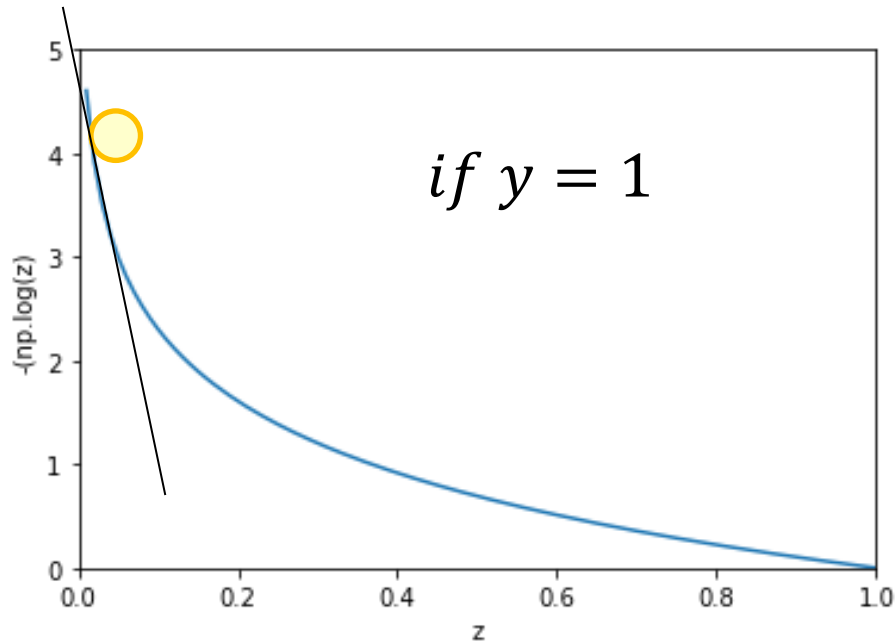
$$J(w) = \frac{1}{2} (\text{sigmoid}(\hat{y}) - y)^2$$

$$\text{Cost}(\hat{y}, y) = \begin{cases} -\log(\hat{y}) & \text{if } y = 1 \\ -\log(1 - \hat{y}) & \text{if } y = 0 \end{cases}$$

$$\text{cost}(\hat{y}, y) = -(y * \log(a) + (1 - y) \log(1 - a))$$

Logistic regression cost function

$$\text{Cost}(\hat{y}, y) = \begin{cases} -\log(\hat{y}) & \text{if } y = 1 \\ -\log(1 - \hat{y}) & \text{if } y = 0 \end{cases}$$

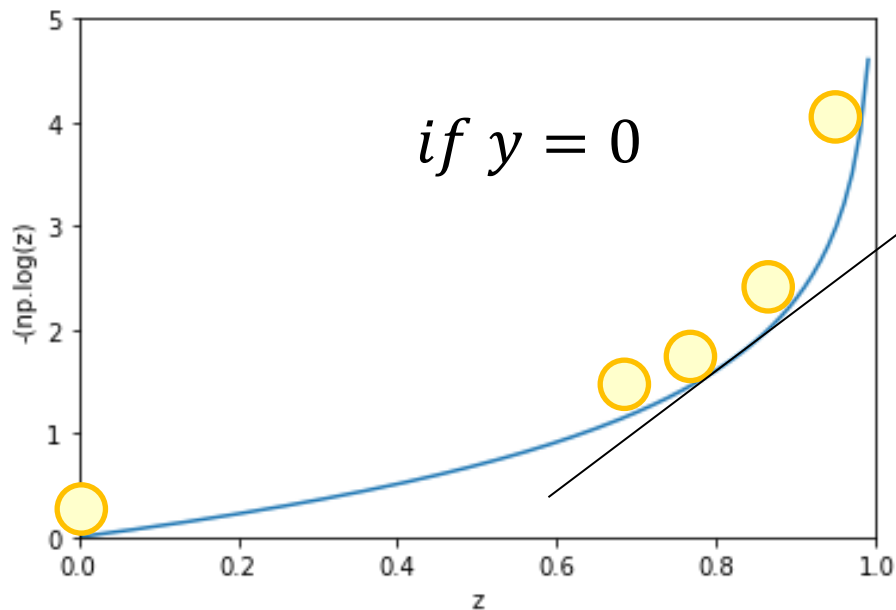


Cost = 0 if $y = 1$, $\hat{y} = 1$

But as $\hat{y} \rightarrow 0$
Cost $\rightarrow \infty$

Logistic regression cost function

$$\text{Cost}(\hat{y}, y) = \begin{cases} -\log(\hat{y}) & \text{if } y = 1 \\ -\log(1 - \hat{y}) & \text{if } y = 0 \end{cases}$$

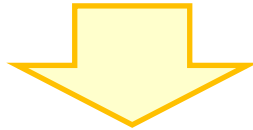


Cost = 0 if $y = 0, \hat{y} = 0$

But as $\hat{y} \rightarrow 1$
Cost $\rightarrow \infty$

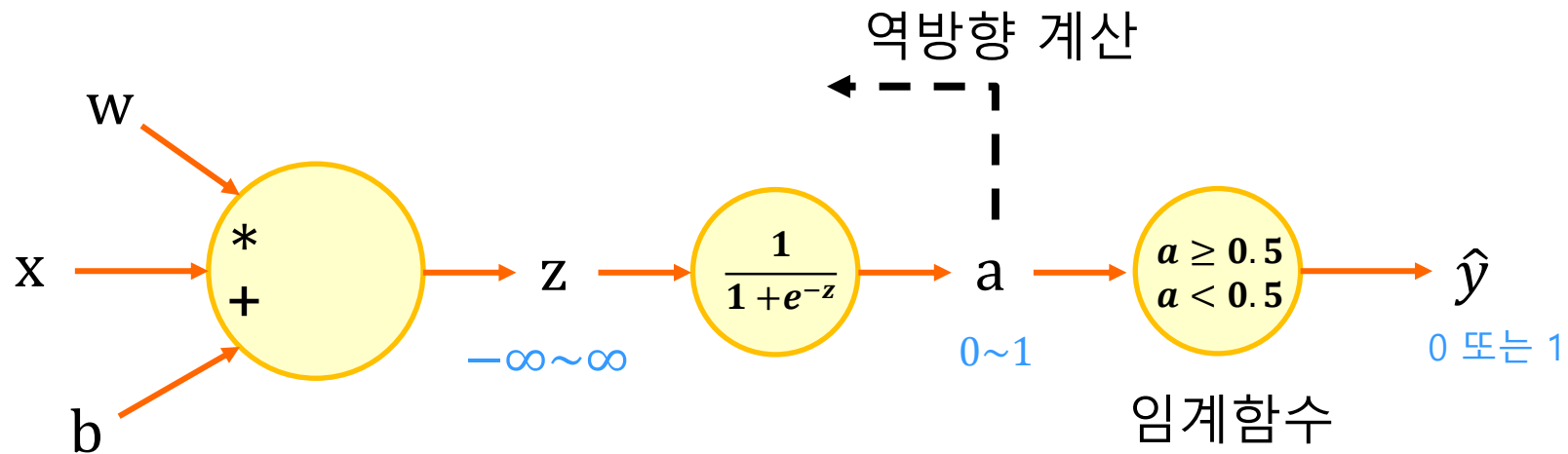
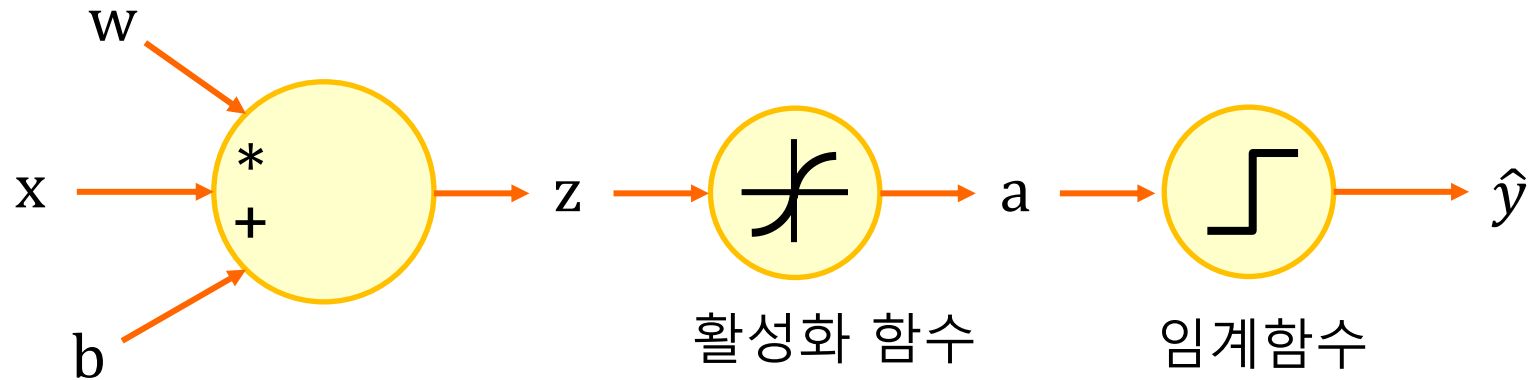
Logistic regression cost function

$$\text{Cost}(\hat{y}, y) = \begin{cases} -\log(\hat{y}) & \text{if } y = 1 \\ -\log(1 - \hat{y}) & \text{if } y = 0 \end{cases}$$



$$L = -(y * \log(a) + (1 - y) \log(1 - a))$$

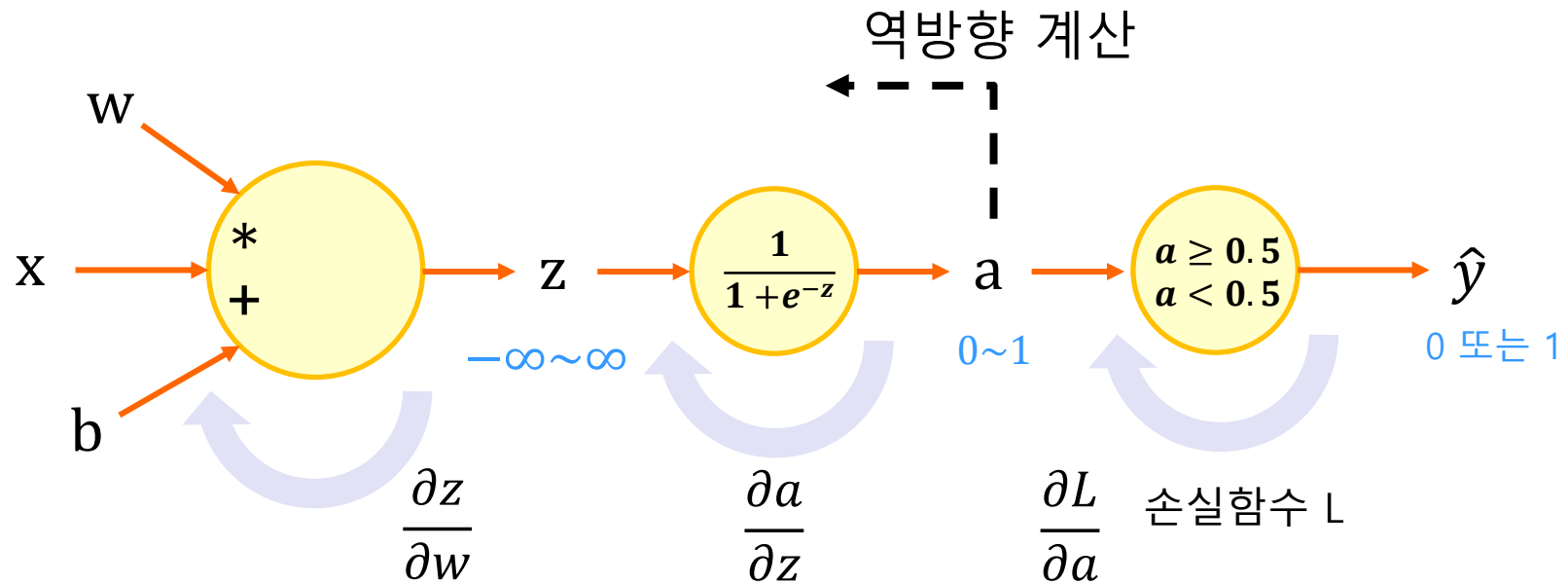
로지스틱 회귀 손실 함수



로지스틱 회귀 손실 함수

특성이 하나인 경우 => x 1개

$$z = wx + b$$



chain rule을 이용하여 각 단계의 미분 결과를 곱한다.

w에 대하여 미분

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial a} \frac{\partial a}{\partial z} \frac{\partial z}{\partial w} = (a - y)x$$

$$\frac{\partial L}{\partial a} = -\left(y \frac{1}{a} - (1 - y) \frac{1}{1-a}\right)$$

$$\frac{\partial a}{\partial z} = a(1 - a)$$

$$\frac{\partial z}{\partial w} = x$$

b에 대하여 미분

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial a} \frac{\partial a}{\partial z} \frac{\partial z}{\partial b} = (a - y)1$$

$$\frac{\partial L}{\partial a} = -\left(y \frac{1}{a} - (1 - y) \frac{1}{1-a}\right)$$

$$\frac{\partial a}{\partial z} = a(1 - a)$$

$$\frac{\partial z}{\partial b} = 1$$

특성이 두개인 경우 => x 2개

$$y = w_1x_1 + w_2x_2 + b$$

$$\frac{\partial}{\partial \hat{y}} \frac{1}{2} (\hat{y} - y)^2 = \hat{y} - y$$

$$\frac{\partial}{\partial w_1} w_1x_1 + w_2x_2 + b = x_1$$

$$\frac{\partial J}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial w_1} = x_1(\hat{y} - y)$$

Chain Rule 사용

$$\hat{y} = w_1x_1 + w_2x_2 + b$$

$$J(w_1, w_2, b) = \frac{1}{2} (\hat{y} - y)^2$$

$$\frac{\partial}{\partial w_1} J(w_1, w_2, b) = \frac{\partial J}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial w_1}$$

특성이 두개인 경우 => x 2개

$$y = w_1x_1 + w_2x_2 + b$$

$$\frac{\partial}{\partial \hat{y}} \frac{1}{2} (\hat{y} - y)^2 = \hat{y} - y$$

$$\frac{\partial}{\partial w_2} w_1x_1 + w_2x_2 + b = x_2$$

$$\frac{\partial J}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial w_2} = x_2(\hat{y} - y)$$

Chain Rule 사용

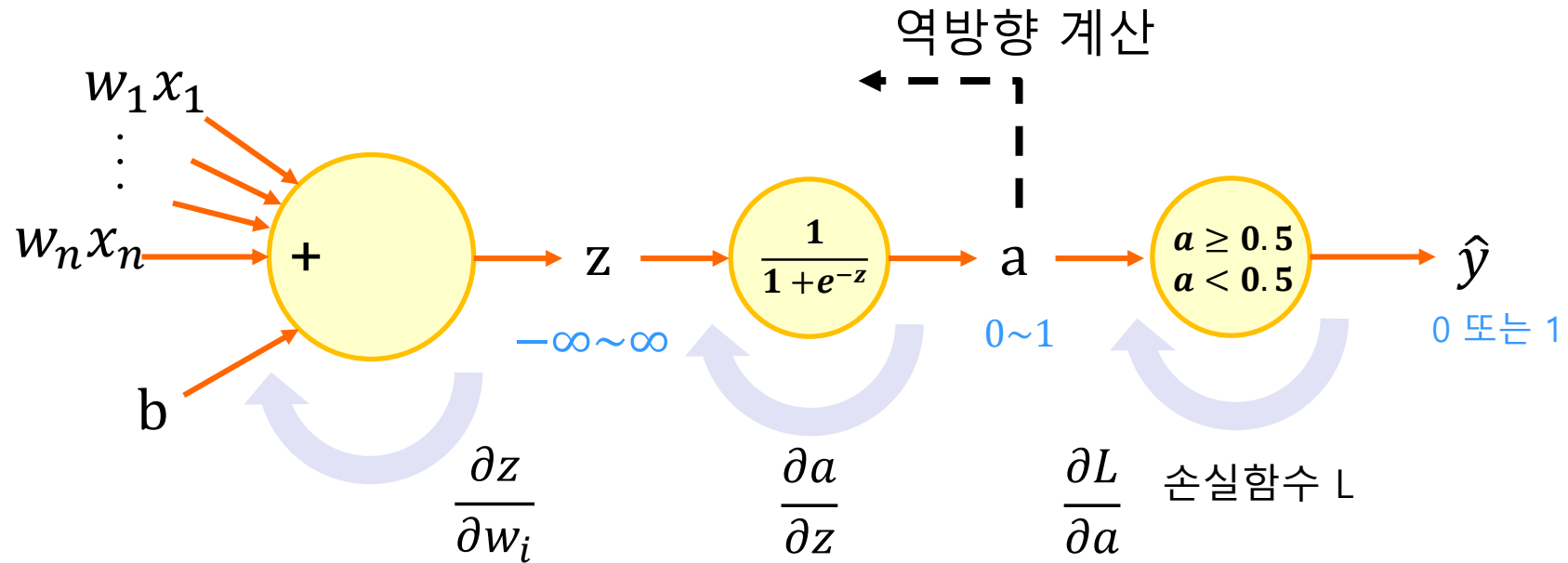
$$\hat{y} = w_1x_1 + w_2x_2 + b$$

$$J(w_1, w_2, b) = \frac{1}{2} (\hat{y} - y)^2$$

$$\frac{\partial}{\partial w_2} J(w_1, w_2, b) = \frac{\partial J}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial w_2}$$

특성이 여러개인 경우 => x n개

$$z = w_1x_1 + w_2x_2 + \dots + w_nx_n + b$$



chain rule을 이용하여 각 단계의 미분 결과를 곱한다.

w에 대하여 미분

$$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial a} \frac{\partial a}{\partial z} \frac{\partial z}{\partial w_i} = (a - y)x_i$$

$$\frac{\partial L}{\partial a} = -\left(y \frac{1}{a} - (1 - y) \frac{1}{1-a}\right)$$

$$\frac{\partial a}{\partial z} = a(1 - a)$$

$$\frac{\partial z}{\partial w_i} = x_i$$

b에 대하여 미분

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial a} \frac{\partial a}{\partial z} \frac{\partial z}{\partial b} = (a - y)1$$

$$\frac{\partial L}{\partial a} = -\left(y \frac{1}{a} - (1 - y) \frac{1}{1-a}\right)$$

$$\frac{\partial a}{\partial z} = a(1 - a)$$

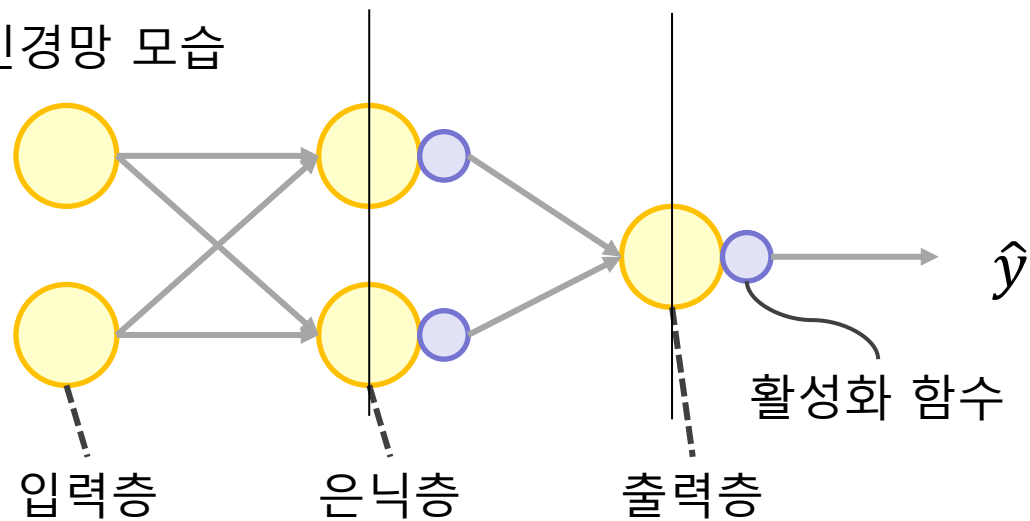
$$\frac{\partial z}{\partial b} = 1$$

제곱 오차의 미분과 로지스틱 손실 함수의 미분은 \hat{y} 이 a 로 바뀌었을 뿐 동일 하다.
따라서 선형함수의 결과 값을 activation 함수인 시그모이드를 적용한 값이 a 이다.

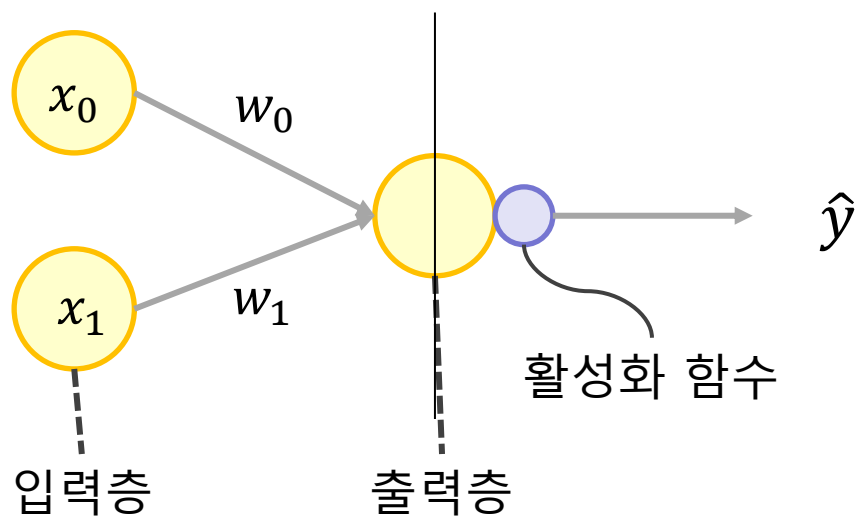
	제곱 오차의 미분	로지스틱 손실 함수의 미분
가중치에 대한 미분	$\frac{\partial SE}{\partial w_i} = (\hat{y} - y)x_i$	$\frac{\partial L}{\partial w_i} = (a - y)x_i$
절편에 대한 미분	$\frac{\partial SE}{\partial b} = (\hat{y} - y)1$	$\frac{\partial L}{\partial b} = (a - y)1$

신경망 구조

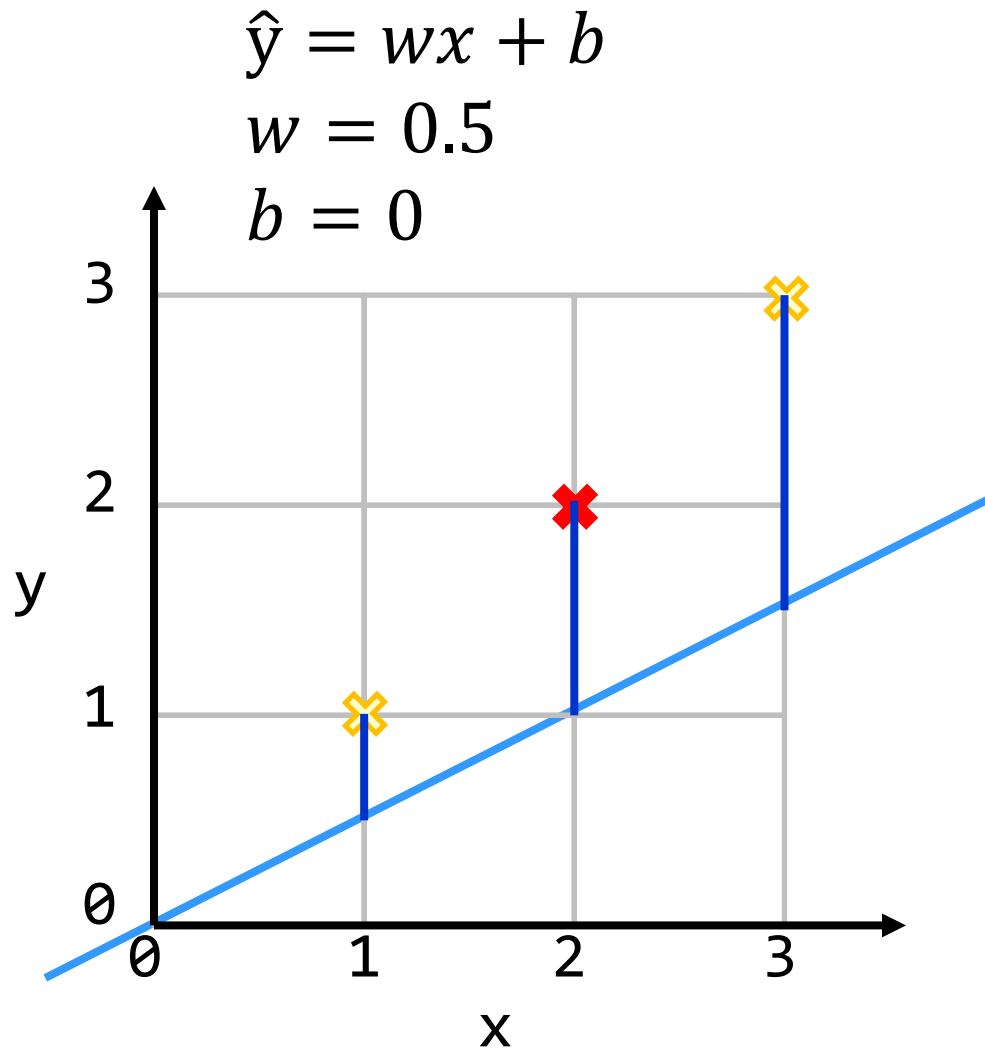
일반적인 신경망 모습



단일층 신경망



신경망 구조



$$\frac{1}{2}(\hat{y} - y)^2 \quad \text{확률적 경사하강}$$

미니배치 경사하강

$$\frac{1}{2m} \sum_{i=1}^m (\hat{y}_i - y_i)^2 \quad \text{배치 경사하강}$$

$$(0.25 + 1 + 2.25)$$

여러 가지 경사 하강법

1번째 샘플 ->	181	92	130	27	...
2번째 샘플 ->	172	56	125	30	...
3번째 샘플 ->	164	61	123	16	...
					...

1개의 샘플을 중복되지 않도록 무작위로 선택 -> 그래디언트 계산

확률적 경사 하강법

1번째 샘플 ->	181	92	130	27	...
2번째 샘플 ->	172	56	125	30	...
3번째 샘플 ->	164	61	123	16	...
					...

전체 샘플들 모두 선택 -> 그래디언트 계산(에포크)

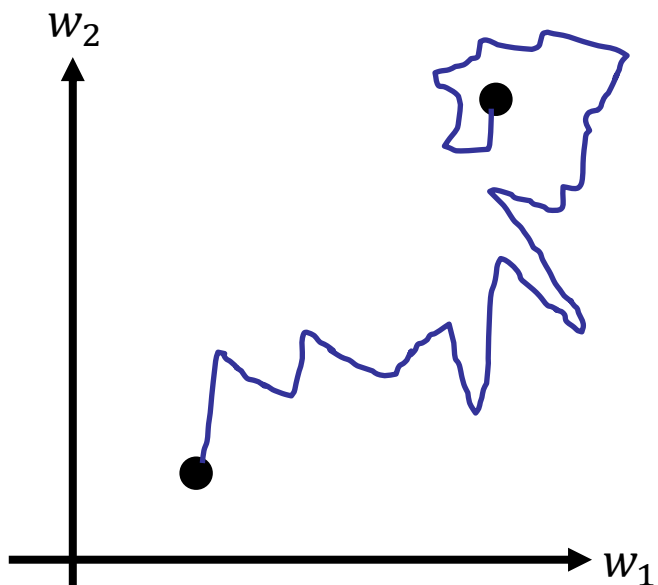
배치 경사 하강법

여러 가지 경사 하강법

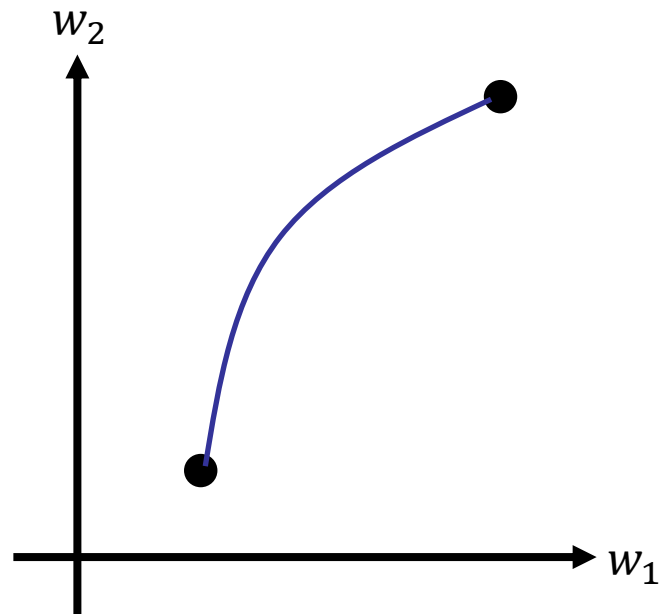
1번째 샘플 ->	181	92	130	27	...
2번째 샘플 ->	172	56	125	30	...
3번째 샘플 ->	164	61	123	16	...
					...

전체 샘플 중 몇 개의 샘플을 중복되지 않도록 > 그레디언트 계산 무작위로 선택

미니 배치 경사 하강법



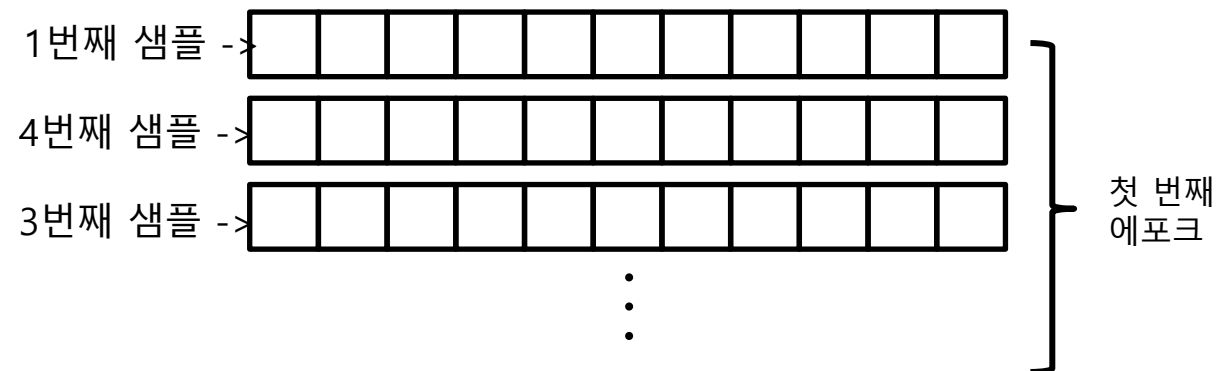
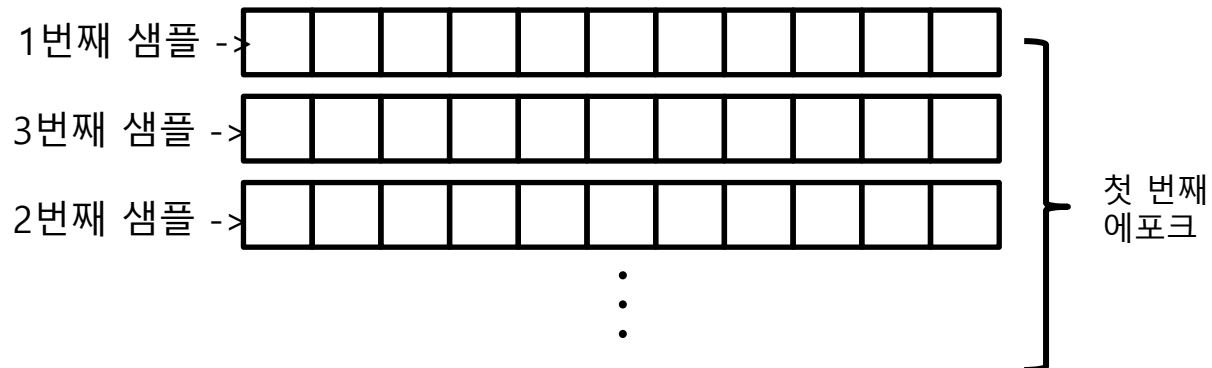
확률적 경사 하강법



배치 경사 하강법

매 에포크마다 훈련 세트의 샘플 순서 섞기

SGD : 확률적 경사하강

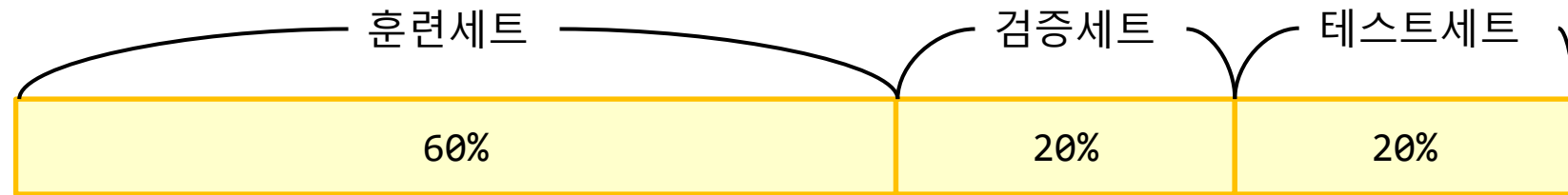


`np.random.permutation()` 함수를 사용하여 인덱스를 섞을수 있다.

검증 세트 분리

"테스트 세트로 모델을 튜닝하면 실전에서 좋은 성능을 기대하기 어렵다"

검증 세트 분리



데이터 전처리와 특성의 스케일

	당도	무게	...
사과1	4	540	...
사과2	8	700	...
사과3	2	480	...

사과의 당도는 1~10이고 사과의 무게의 범위는 500~1000이다.
이런 경우 '두 특성의 스케일 차이가 크다'라고 말한다.

스케일 조정

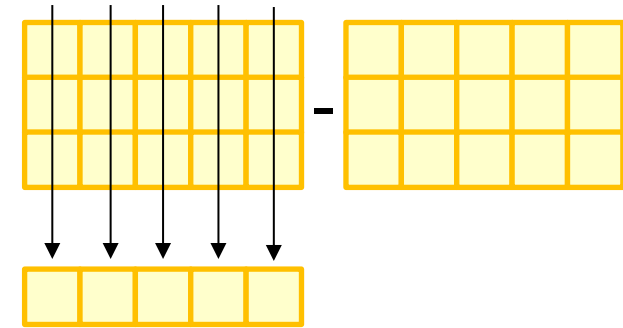
$$Z = \frac{X - \mu}{s}$$

표준화는 특성 값에서 평균을 빼고
표준 편차로 나누면 된다.

$$s = \sqrt{\frac{1}{m} \sum_{i=0}^m (x_i - \mu)^2}$$

표준 편차 공식

(3, 5)



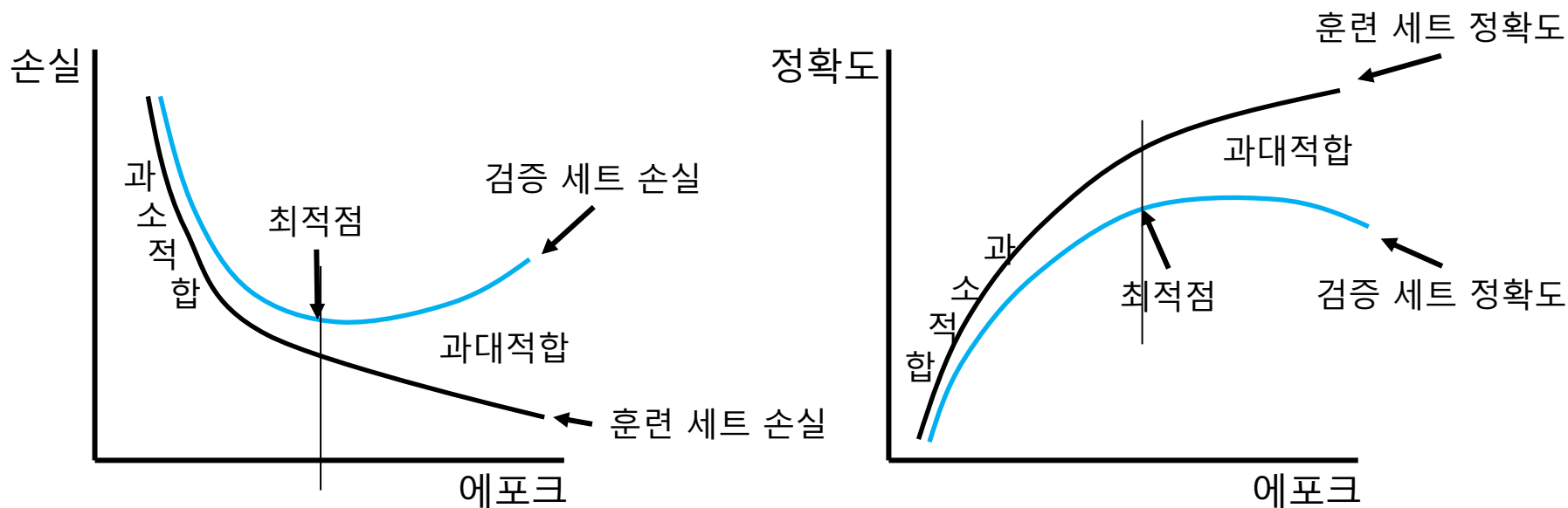
```
train_mean = np.mean(x_train, axis=0)
train_std = np.std(x_train, axis=0)
x_train_scaled = (x_train - train_mean) / train_std
```


과대적합 :

훈련 세트에서는 좋은 성능을 내지만 검증 세트에서는 낮은 성능을 내는 경우

과소적합 :

훈련 세트와 검증세트의 성능에는 차이가 크지 않지만 모두 낮은 성능을 내는 경우



L1 규제는 손실 함수에 가중치의 절대값인 L1 노름(norm)을 추가한다.

$$\|w\|_1 = \sum_{i=1}^n |w_i|$$

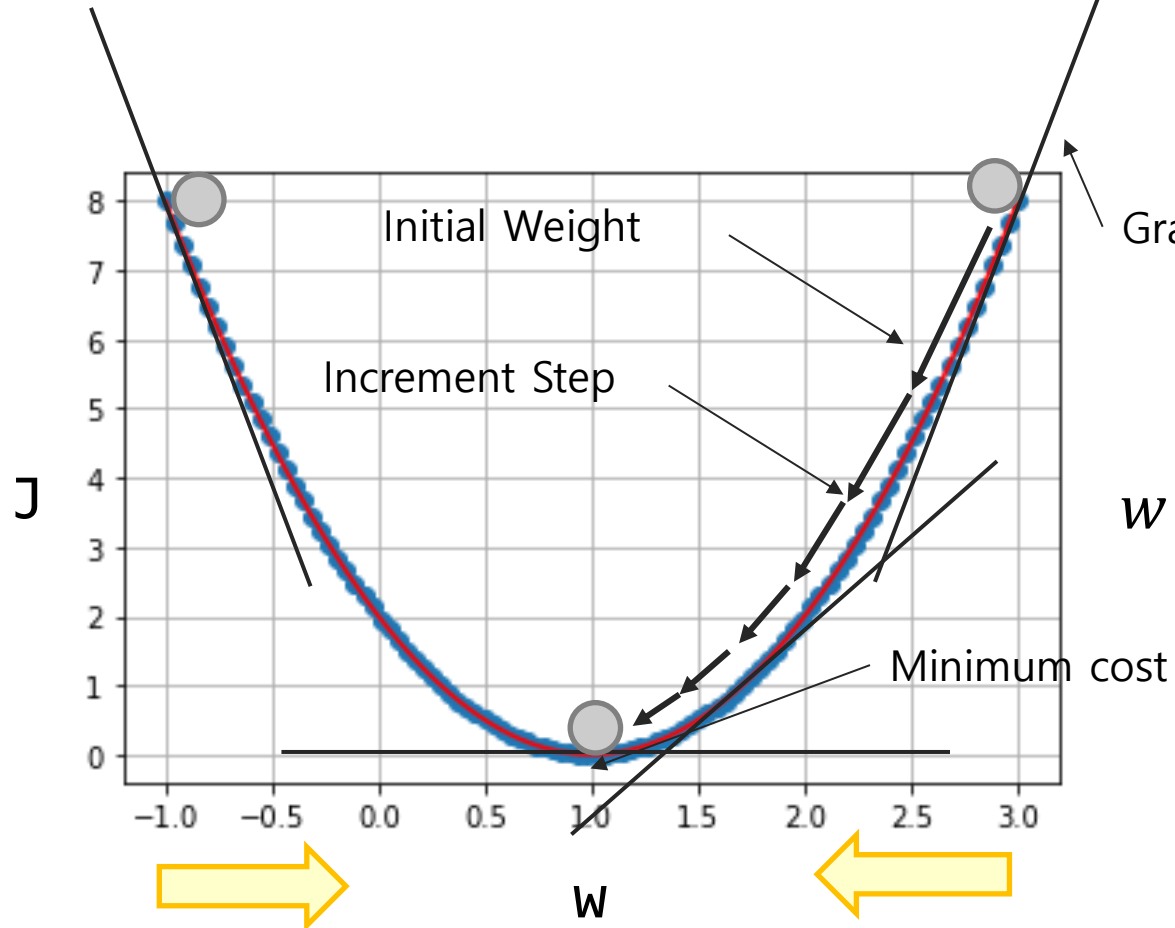
$$L = -(y \log(a) + (1 - y) \log(1 - a))$$

L1 노름을 그냥 더하지 않고 규제의 양을 조절하는 파라미터 α 를 곱한 후 더한다.

$$L = -(y \log(a) + (1 - y) \log(1 - a)) + \alpha \sum_{i=1}^n |w_i|$$

Gradient descent algorithm

$$\hat{y} = wx + b$$



$$w = w - \alpha(\hat{y} - y)x$$

$$0.08$$

$$w = w - \alpha(\hat{y} - y)x + 0.01 * \text{sign}(w)$$

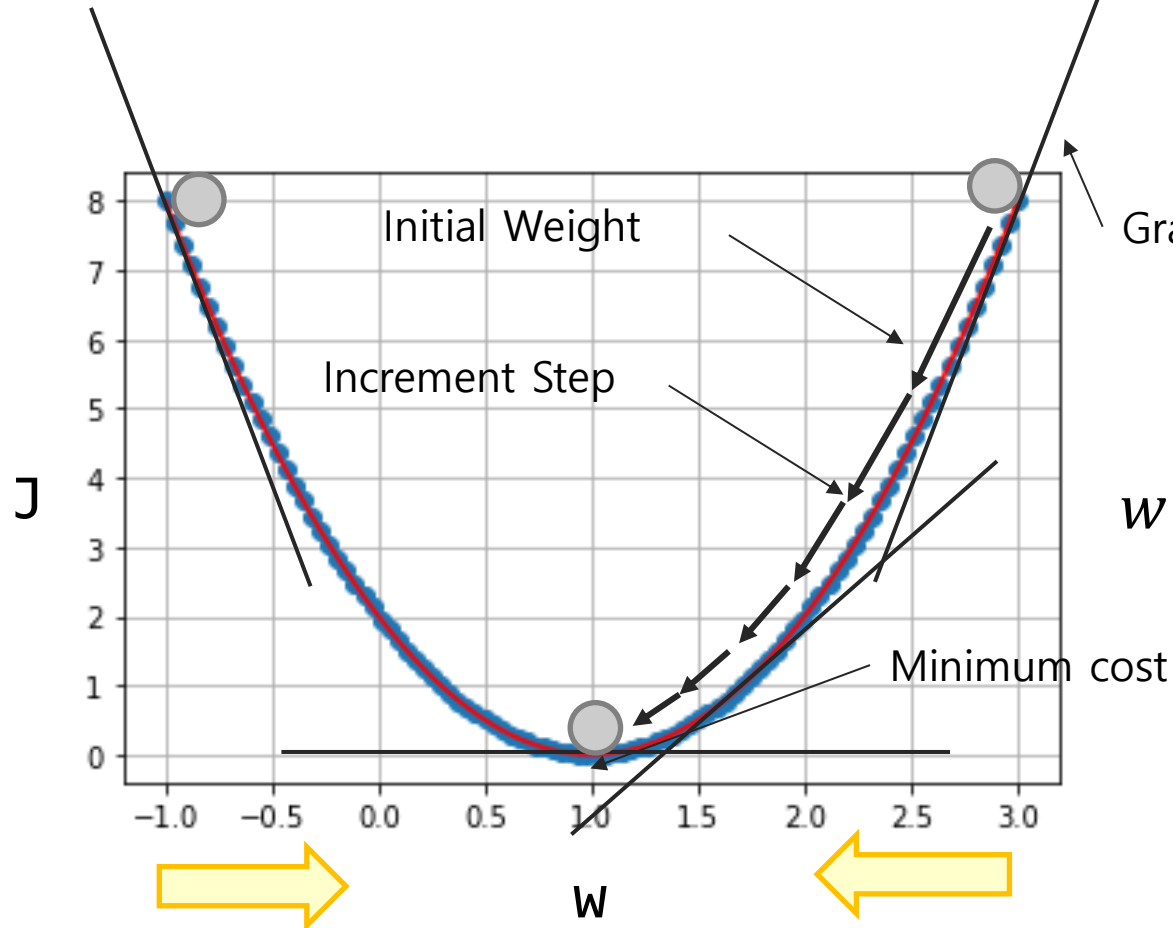
$$-0.08 + 0.01$$

$$x = 2$$

$$y = 2$$

Gradient descent algorithm

$$\hat{y} = wx + b$$



$$w = w - \alpha(\hat{y} - y)x$$

$$0.08$$

$$w = w - \alpha(\hat{y} - y)x + 0.01 * \text{sign}(w)$$

$$+0.08 - 0.01$$

$$x = 2$$

$$y = 2$$

L1 규제는 손실 함수에 가중치의 절대값인 L1 노름(norm)을 추가한다.

$$\frac{\partial}{\partial w} L = (a - y)x + \alpha * \text{sign}(w)$$

$$w = w - \eta \frac{\partial L}{\partial w} = w - \eta(a - y)x + \alpha * \text{sign}(w)$$

파이썬으로 작성된 L1 규제 적용된 오차 역전파 구현

```
w_grad += alpha * np.sign(w)
```

회귀 모델에 L1 규제를 추가한 것을 라쏘 모델이라 한다.

L2 규제는 손실 함수에 가중치에 대한 L2 노름(norm)의 제곱을 더한다.

$$\|w\|_2 = \sum_{i=1}^n |w_i|^2$$

$$L = -(y \log(a) + (1 - y) \log(1 - a))$$

L1 노름을 그냥 더하지 않고 규제의 양을 조절하는 파라미터 α 를 곱한 후 더한다.

$$L = -(y \log(a) + (1 - y) \log(1 - a)) + \frac{1}{2} \alpha \sum_{i=1}^n |w_i|^2$$

L2 규제는 손실 함수에 가중치에 대한 L2 노름(norm)의 제곱을 더한다.

$$\frac{\partial}{\partial w} L = -(y - a)x + \alpha * w$$

$$w = w - \eta \frac{\partial L}{\partial w} = w + \eta((y - a)x - \alpha * w)$$

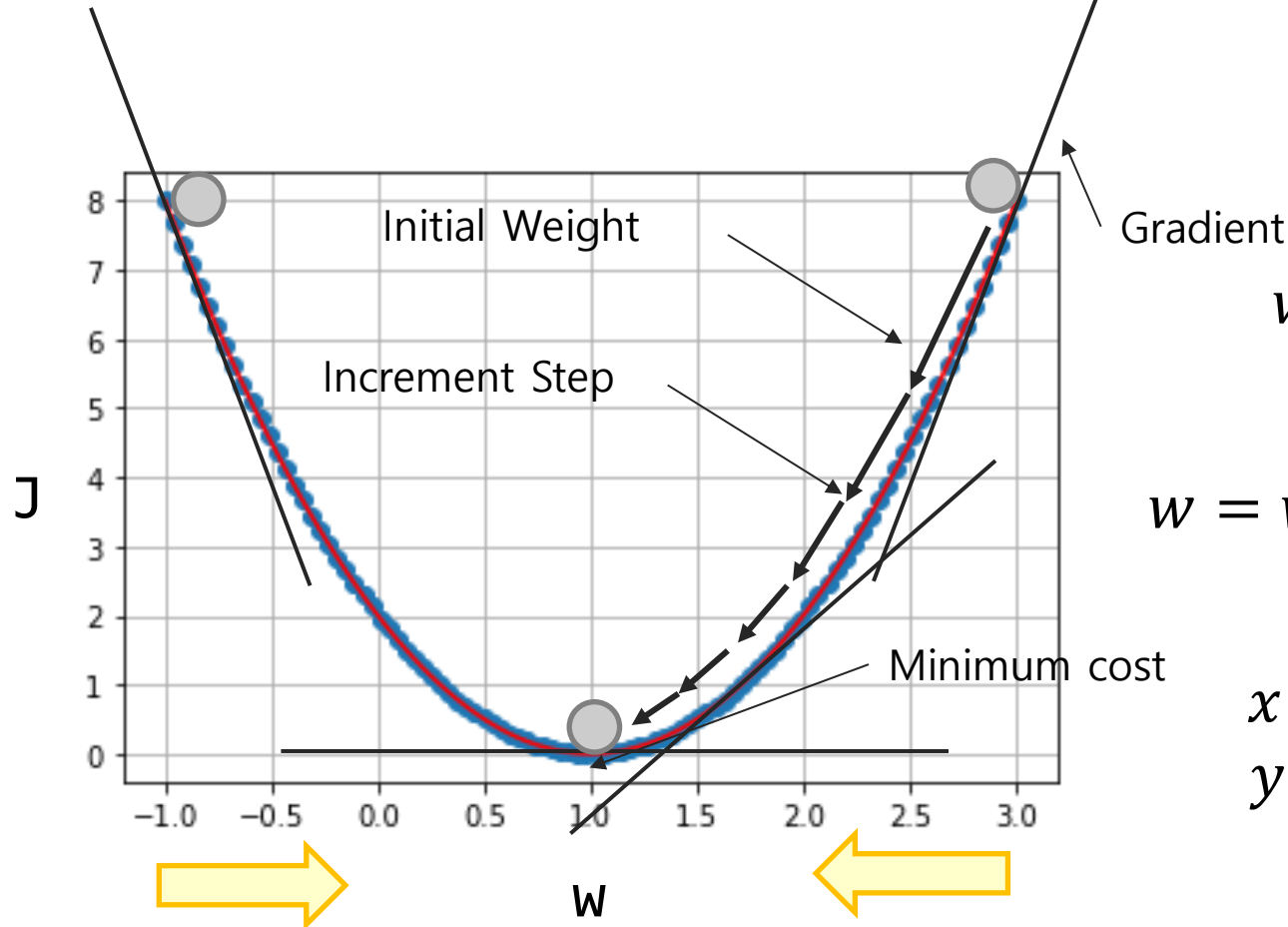
파이썬으로 작성된 L2 규제 적용된 오차 역전파 구현

```
w_grad += alpha * w
```

회귀 모델에 L2 규제를 추가한 것을 릿지 모델이라 한다.

Gradient descent algorithm

$$\hat{y} = wx + b$$



$$w = w - \alpha(\hat{y} - y)x$$

$$0.08$$

$$w = w - \alpha(\hat{y} - y)x + 0.01 * w$$

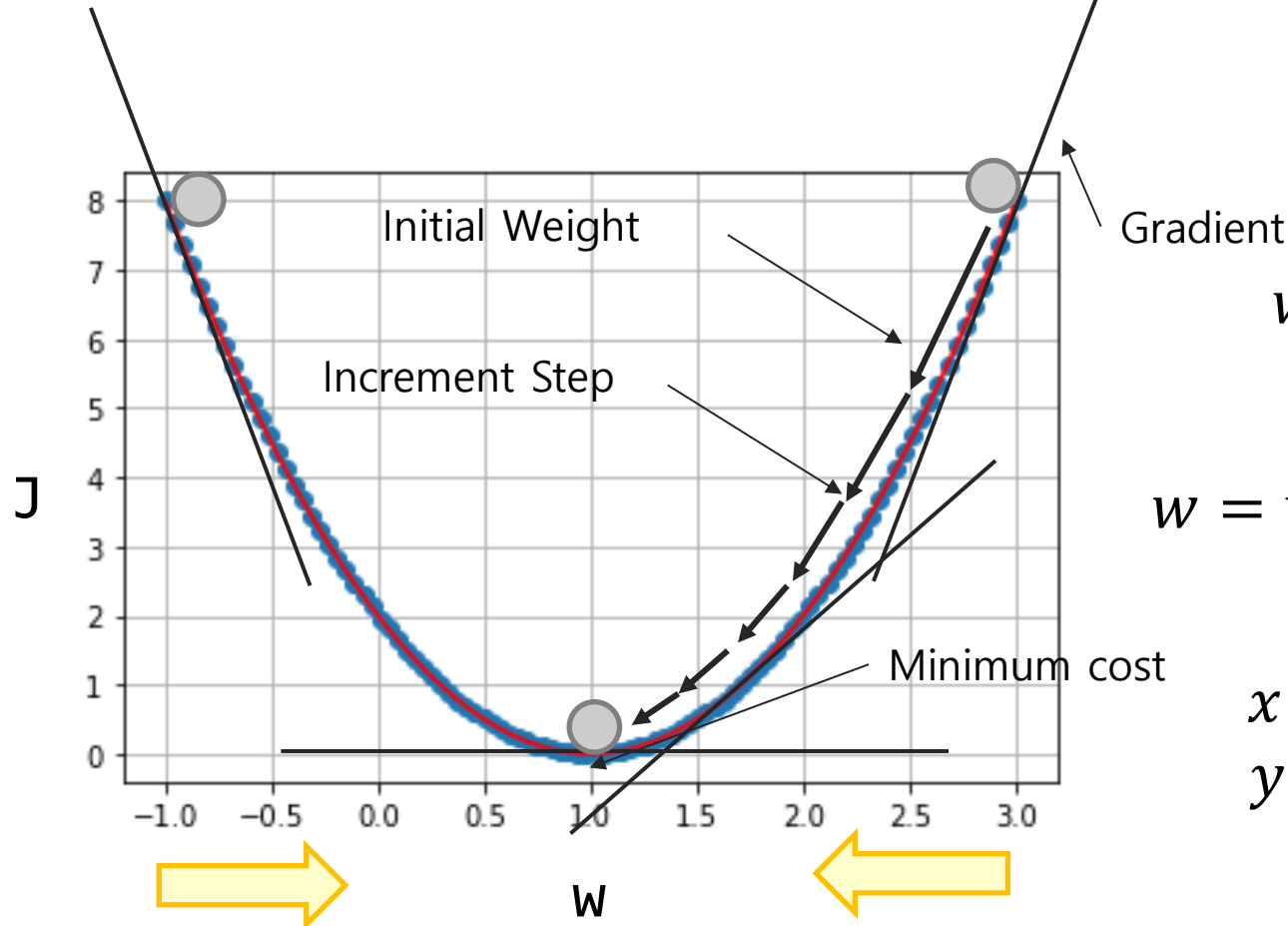
$$-0.08 + 0.03$$

$$x = 2$$

$$y = 2$$

Gradient descent algorithm

$$\hat{y} = wx + b$$



$$w = w - \alpha(\hat{y} - y)x$$

$$0.08$$

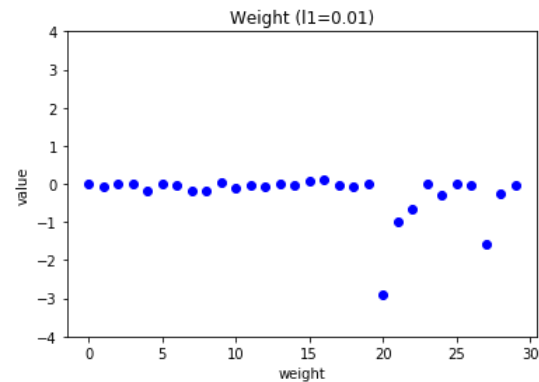
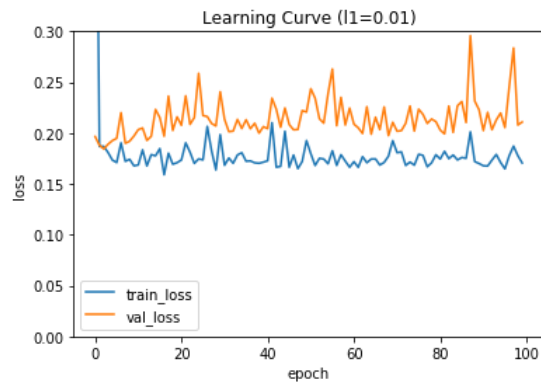
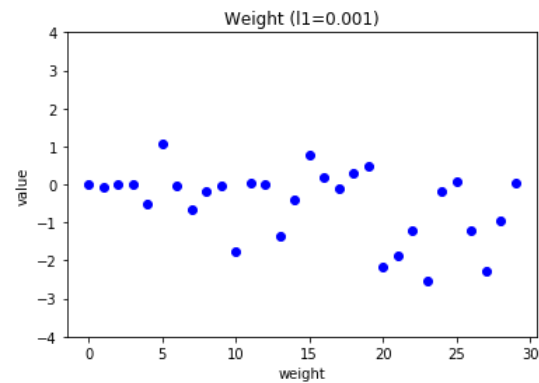
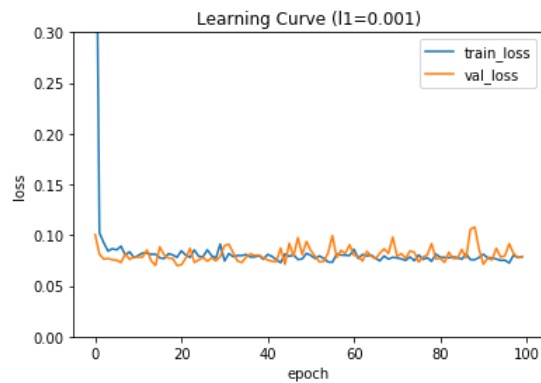
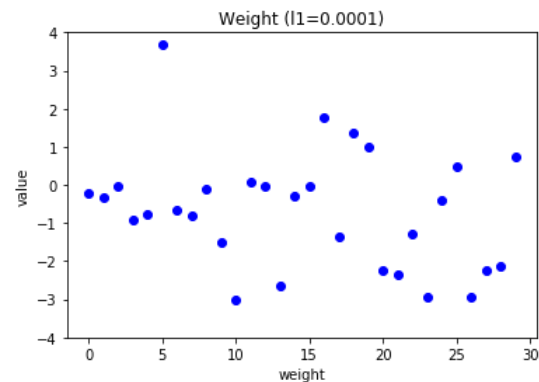
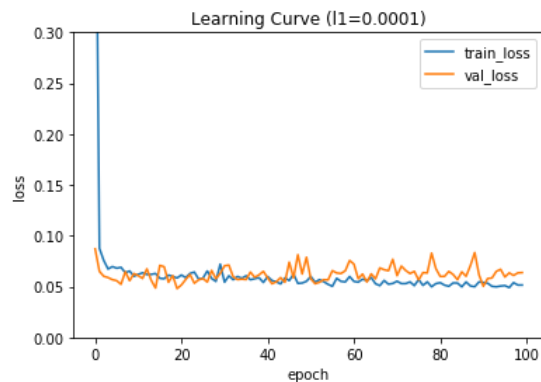
$$w = w - \alpha(\hat{y} - y)x + 0.01 * w$$

$$+0.08 - 0.01$$

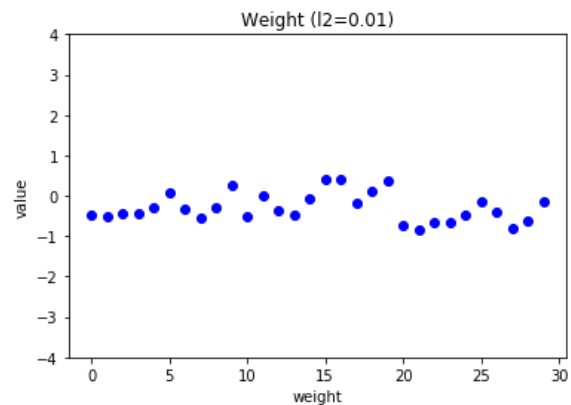
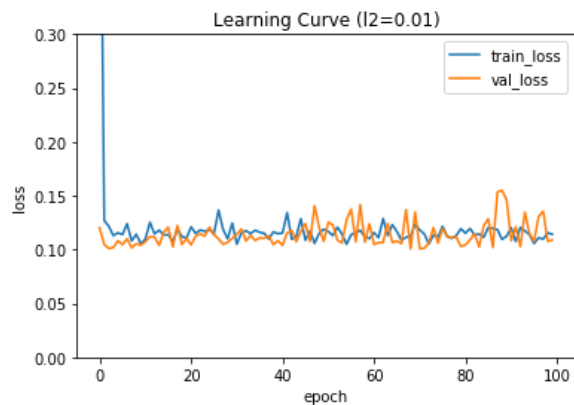
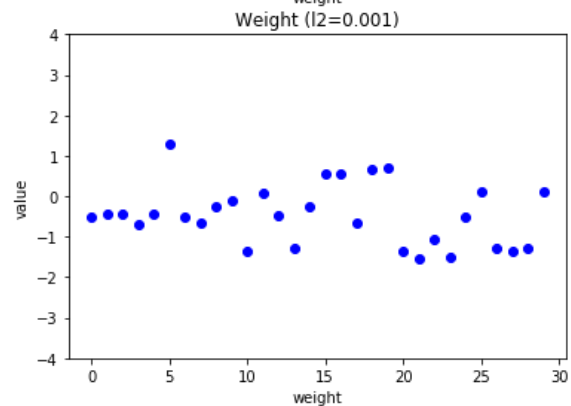
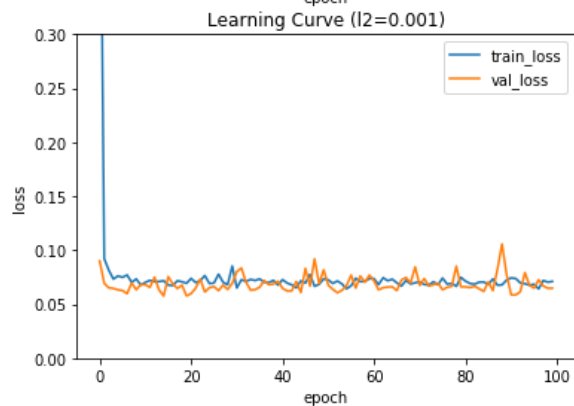
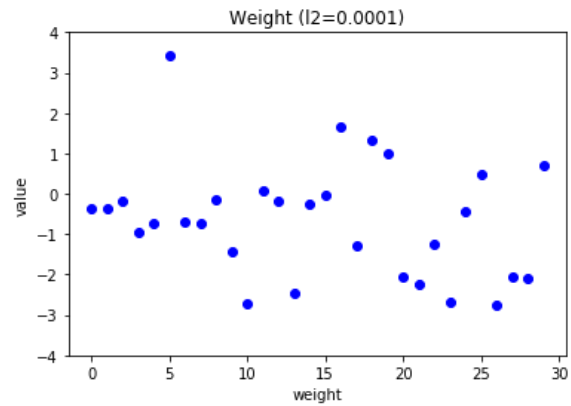
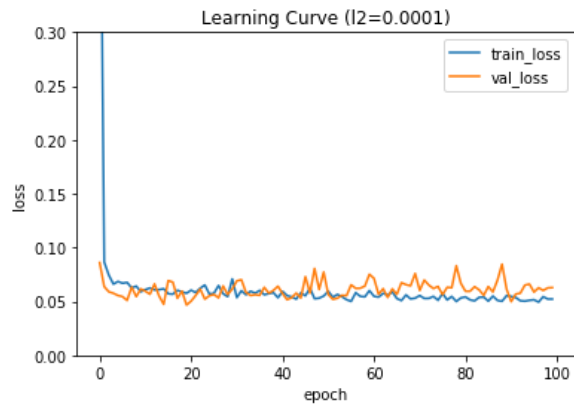
$$x = 2$$

$$y = 2$$

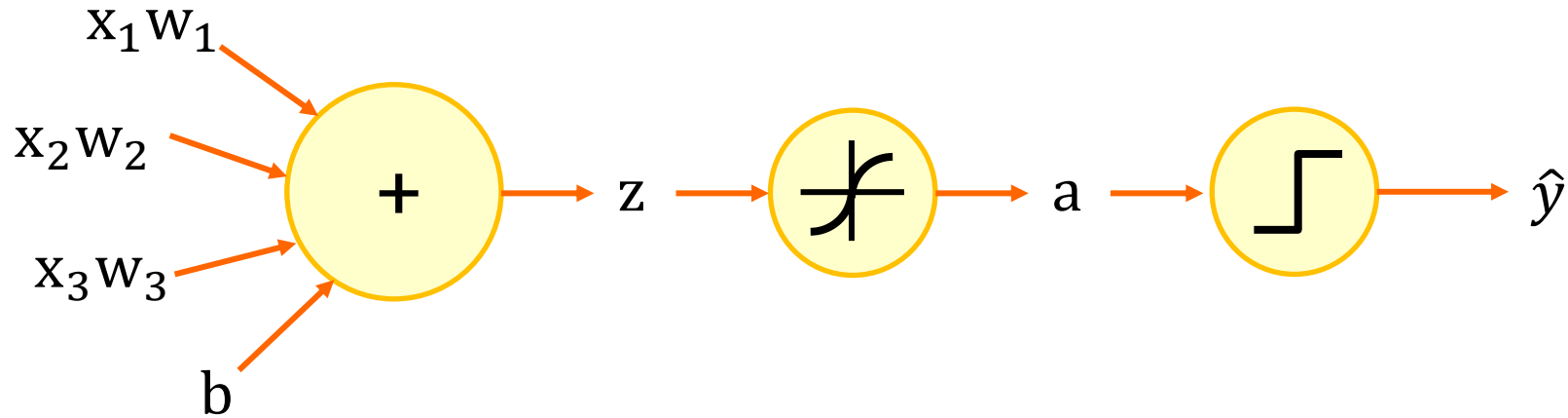
L1 규제 적용



L2 규제 적용



벡터 연산과 행렬 연산

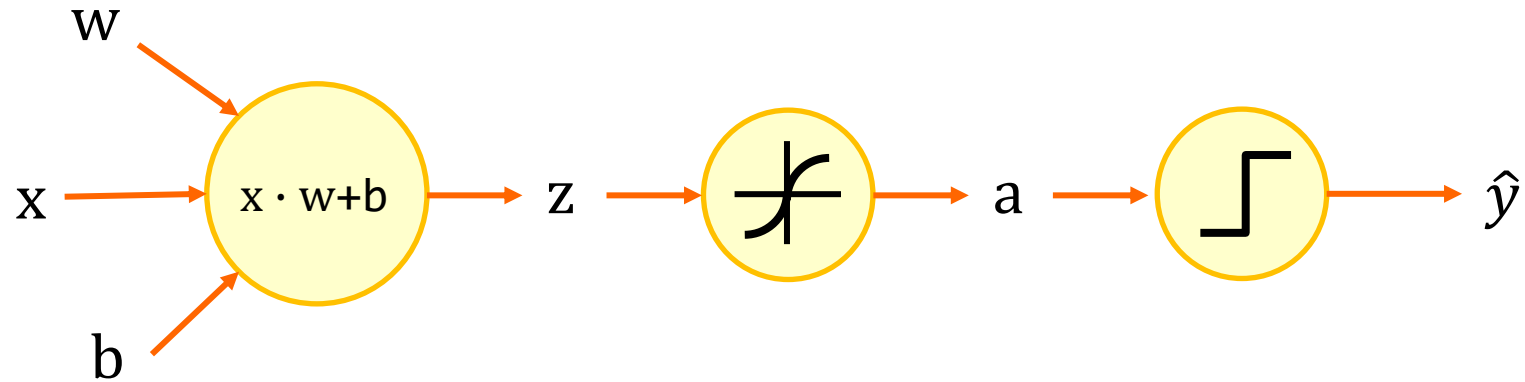


```
def forpass(self, x):  
    z = np.sum(x * self.w) + self.b  
    return z
```

넘파이의 원소별 곱셈

```
x = [x1, x1, ..., xn]  
w = [w1, w, ..., wn]  
x * w = [x1 * w1, x2 * w2, ..., xn * w1]
```

벡터 연산과 행렬 연산



점 곱을 행렬 곱셈으로 표현

$$XW = \begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} = x_1 \times w_1 + x_2 \times w_2 + x_3 \times w_3$$

$(1, 3)(3, 1) \Rightarrow (1, 1)$

```
z = np.dot(x , self.w) + self.b
```

$$XW = \begin{bmatrix} x_1^{(1)} & x_2^{(1)} & x_3^{(1)} \\ x_1^{(2)} & x_2^{(2)} & x_3^{(2)} \\ \vdots & \vdots & \vdots \\ x_1^{(m)} & x_2^{(m)} & x_3^{(m)} \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} = \begin{bmatrix} x_1^{(1)}w_1 + x_2^{(1)}w_2 + x_3^{(1)}w_3 \\ x_1^{(2)}w_1 + x_2^{(2)}w_2 + x_3^{(2)}w_3 \\ \vdots \\ x_1^{(m)}w_1 + x_2^{(m)}w_2 + x_3^{(m)}w_3 \end{bmatrix}$$

행렬곱 가능과 곱 결과 크기

$$(m, n) \cdot (n, k) = (m, k)$$

첫 번째 행렬의 열(n)과 두 번째 행렬의 행(n)의 크기는 반드시 같아야 한다.

곱 결과의 크기는 첫 번째 행렬의 행(m)과 두 번째 행렬의 열(k)의 크기가 된다.

정방향 계산을 행렬 곱셈으로
표현

$$XW = \begin{bmatrix} x_1^{(1)} & \cdots & x_{30}^{(1)} \\ \vdots & \ddots & \vdots \\ x_1^{(364)} & \cdots & x_{30}^{(364)} \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_{30} \end{bmatrix} + \begin{bmatrix} b \\ b \\ \vdots \\ b \end{bmatrix} = \begin{bmatrix} z^{(1)} \\ z^{(2)} \\ \vdots \\ z^{(364)} \end{bmatrix}$$

$$= \begin{bmatrix} x_1^{(1)} & \cdots & x_{30}^{(1)} \\ \vdots & \ddots & \vdots \\ x_1^{(364)} & \cdots & x_{30}^{(364)} \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_{30} \end{bmatrix} + b = \begin{bmatrix} x_1^{(1)}w_1 + x_2^{(1)}w_2 + \cdots + x_{30}^{(1)}w_{30} + b \\ \vdots \\ x_1^{(364)}w_1 + x_2^{(364)}w_2 + \cdots + x_{30}^{(364)}w_{30} + b \end{bmatrix}$$

그레디언트 계산

$$X^T E = \begin{bmatrix} x_1^{(1)} & \cdots & x_1^{(364)} \\ x_2^{(1)} & \cdots & x_2^{(364)} \\ \vdots & \cdots & \vdots \\ x_{30}^{(1)} & \cdots & x_{30}^{(364)} \end{bmatrix} \begin{bmatrix} e^{(1)} \\ e^{(3)} \\ \vdots \\ e^{(364)} \end{bmatrix} = \begin{bmatrix} g_1 \\ g_2 \\ \vdots \\ g_{30} \end{bmatrix}$$

그레디언트는 X 와 E 의 행렬곱이다.
그러나 벡터 연산에서 X 의 크기는 $(364, 30)$ 이고
 E 는 $(364, 1)$ 이므로 행렬의 곱을 할 수 없다.

이때 X 를 전치하면 행과 열이 바뀌므로
 $X^T (30, 364)$ 가 되므로

$X^T E$ 는 $(30, 364) \cdot (364, 1)$ 이므로 곱이 가능하고
결과는 $(30, 1)$ 이 되므로 그레디언트와 같은 행렬을 구할 수 있다.

$$W = W - X^T E$$

$(30, 364)(364, 1)$
 $(30, 1)$

$$X \Rightarrow (364, 30)$$

$$W \Rightarrow (30, \textcolor{red}{1})$$

$$\hat{y} \Rightarrow (364, 1)$$

$$a \Rightarrow (364, 1)$$

$$err \Rightarrow (a - y) \Rightarrow (364, 1)$$

행렬곱의 미분(MatMul)

$$\frac{\partial}{\partial \hat{y}} \frac{1}{2} (\hat{y} - y)^2$$

$$= \hat{y} - y = E$$

$$\frac{\partial}{\partial w} wx + b$$

$$= X$$

$$\frac{\partial J}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial w} = X^T E$$

Chain Rule 사용

$$\hat{y} = XW + b$$

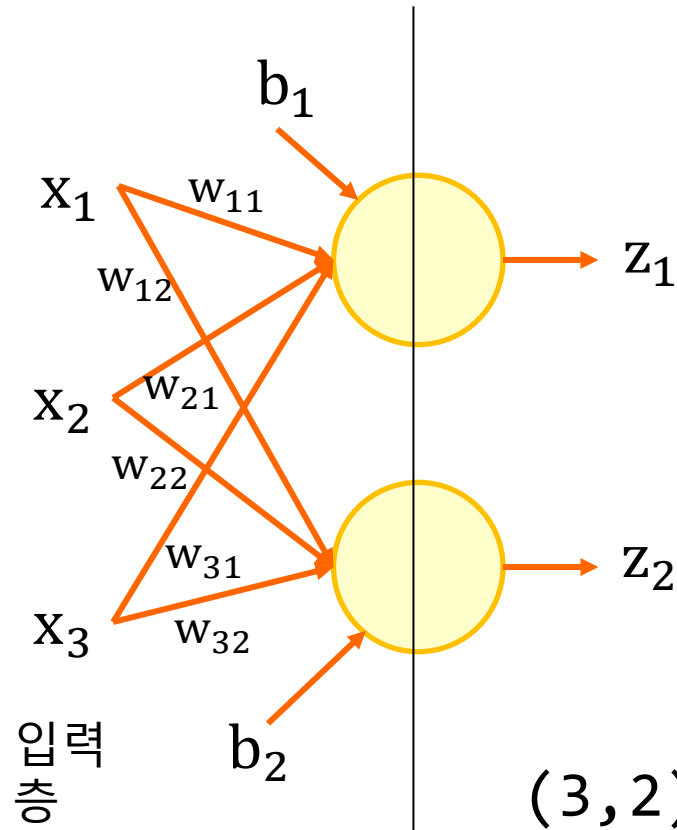
$$J(w,b) = \frac{1}{2} (\hat{y} - y)^2$$

$$\frac{\partial}{\partial w} J(w,b) = \frac{\partial J}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial w}$$

전체식을 w 로 미분한 경우
사라진 w 의 자리에는 뒤에서 날라온
미분값인 E 를 쓰고
남아 있는 x 는 전치하여
행렬 곱을 한다.

<https://playground.tensorflow.org/>

하나의 층에 여러개의 뉴런 사용



$$XW_1 + b_1 = z_1$$

$$XW_2 + b_2 = z_2$$

$(3, 2) \Rightarrow (\text{특성수}, \text{뉴런수}) \Rightarrow (\text{입력수}, \text{출력수})$

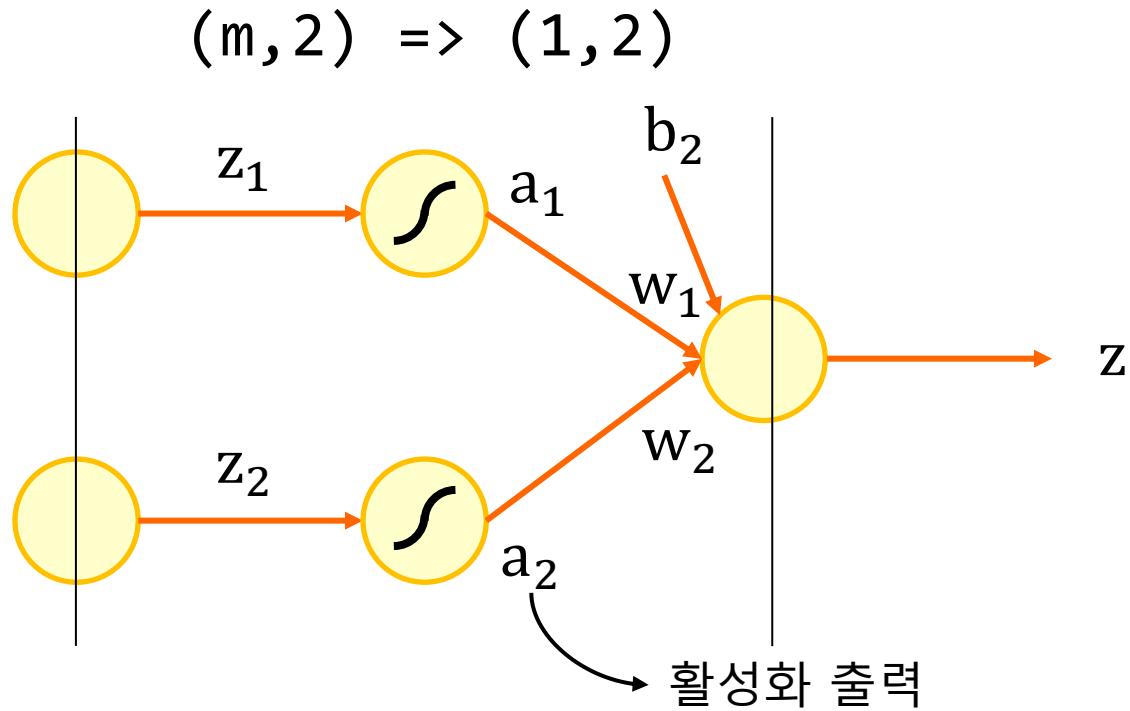
$$x_1 w_{11} + x_2 w_{21} + x_3 w_{31} + b_1 = z_1$$

$$x_1 w_{12} + x_2 w_{22} + x_3 w_{32} + b_2 = z_2$$

$$\begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix} \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \\ w_{31} & w_{32} \end{bmatrix} + \begin{bmatrix} b_1 & b_2 \end{bmatrix} = \begin{bmatrix} z_1 & z_2 \end{bmatrix}$$

$$(m, n) \quad (n, 2) \Rightarrow (m, 2)$$

하나의 층에 여러개의 뉴런 사용



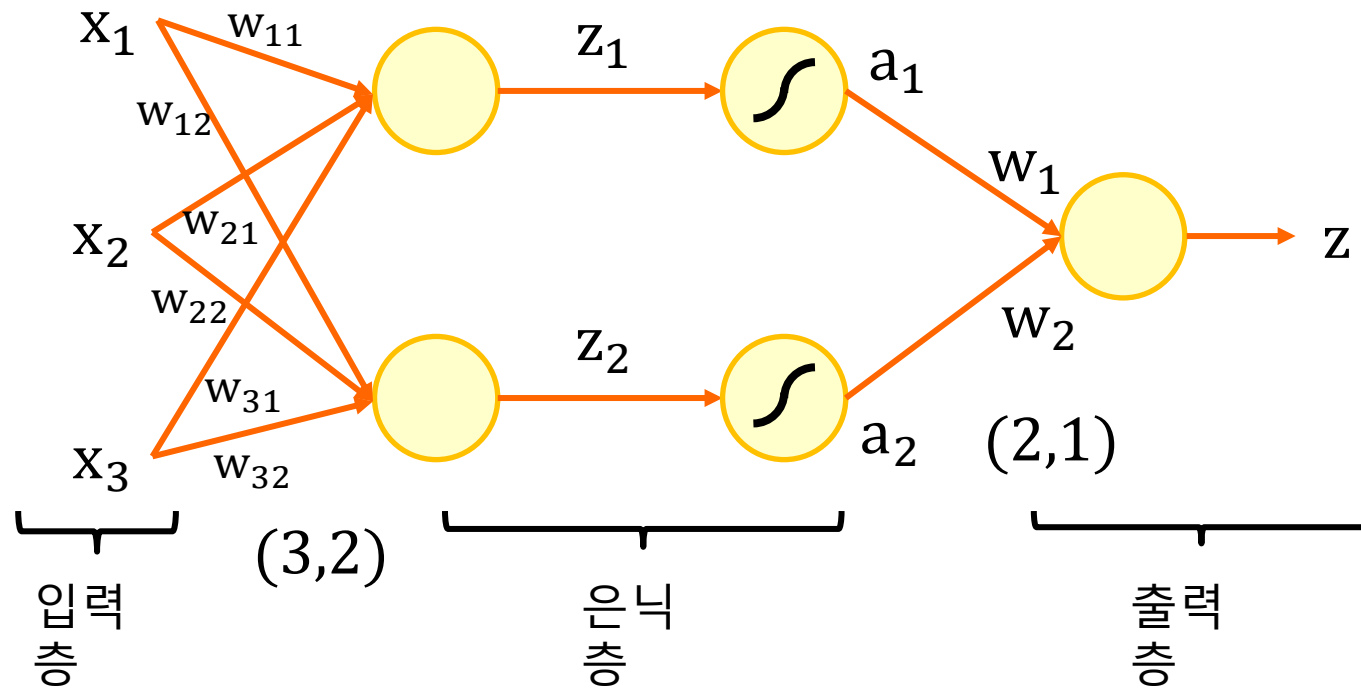
$$a_1 w_1 + a_2 w_2 + b_2 = z$$

$$\begin{bmatrix} a_1 & a_2 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} + b_2 = z$$

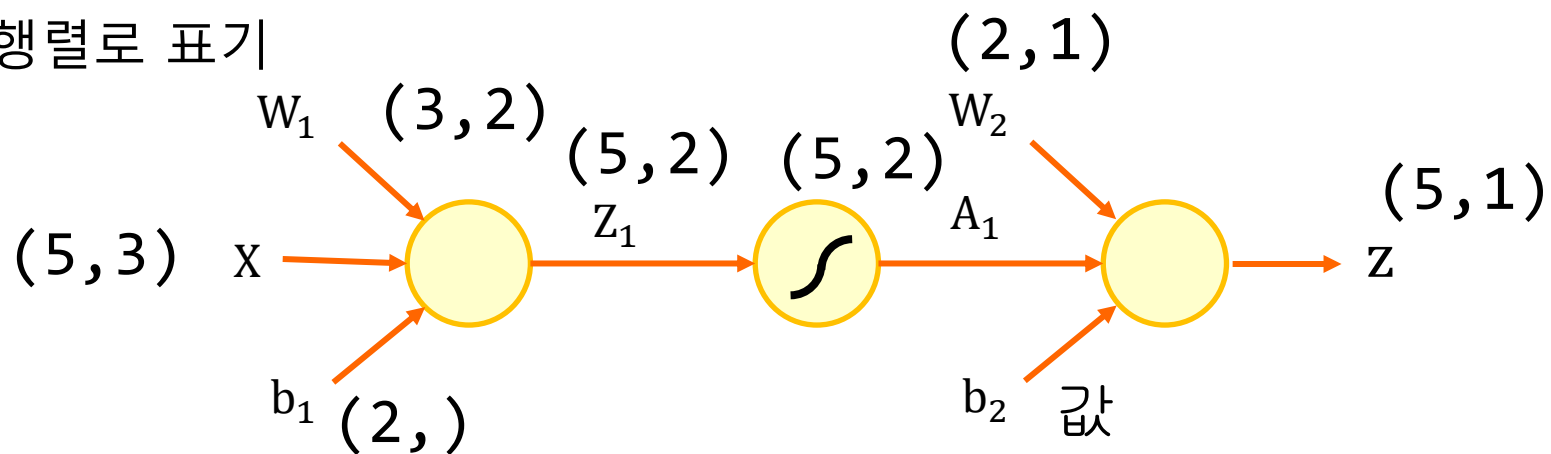
$$(m, 2)(2, 1) \Rightarrow (m, 1)$$

$$A_1 W_2 + b_2 = z_2$$

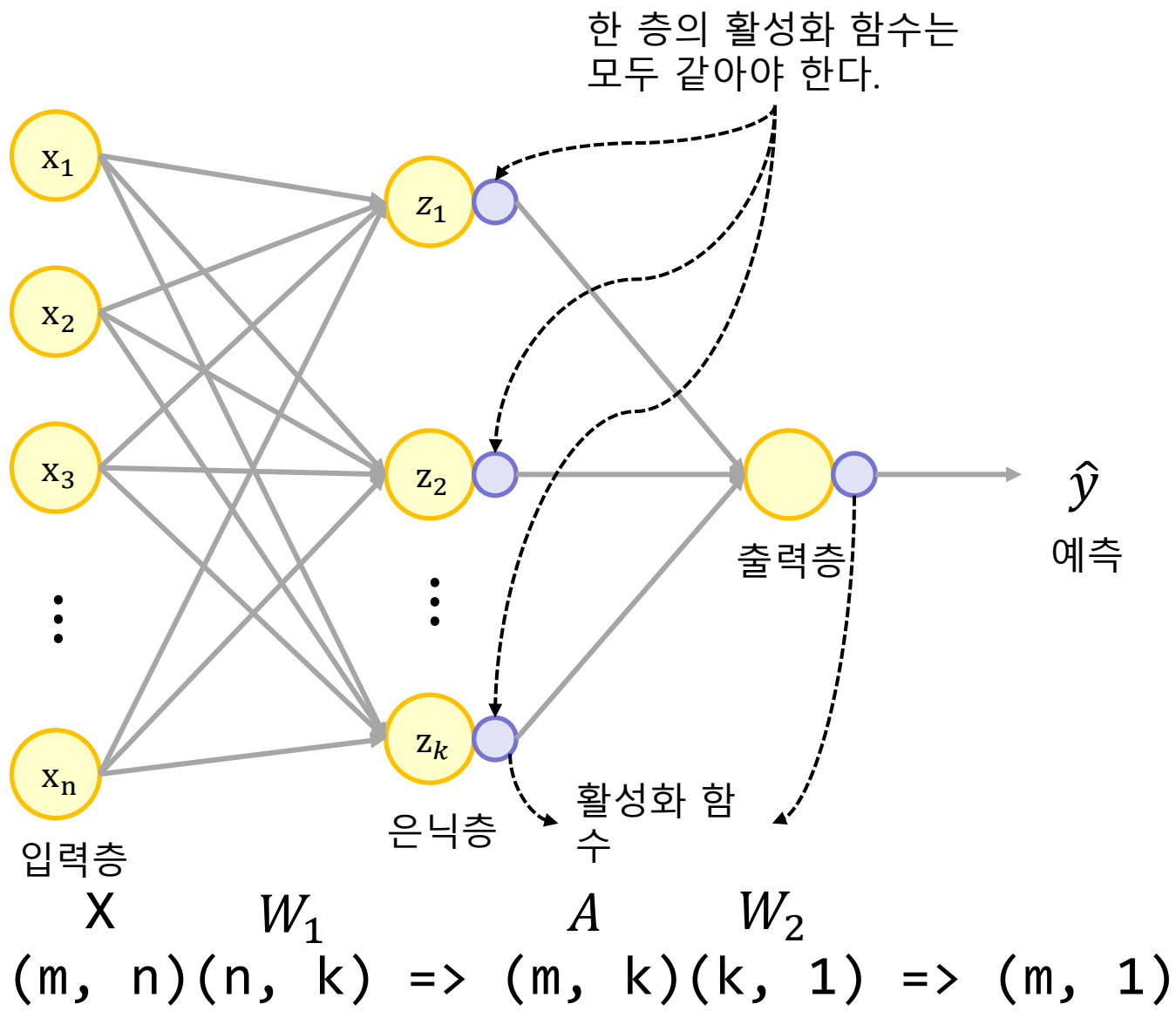
하나의 층에 여러개의 뉴런 사용



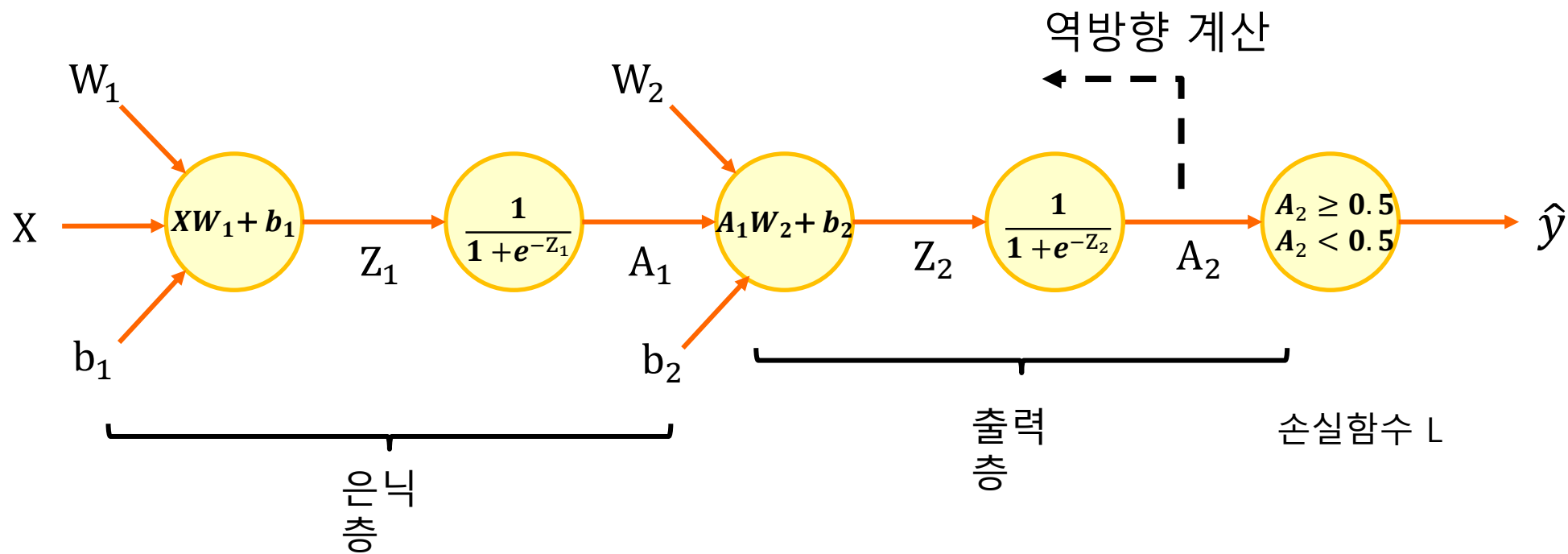
행렬로 표기



2층의 완전연결 신경망



다층 신경망 backprop 이해



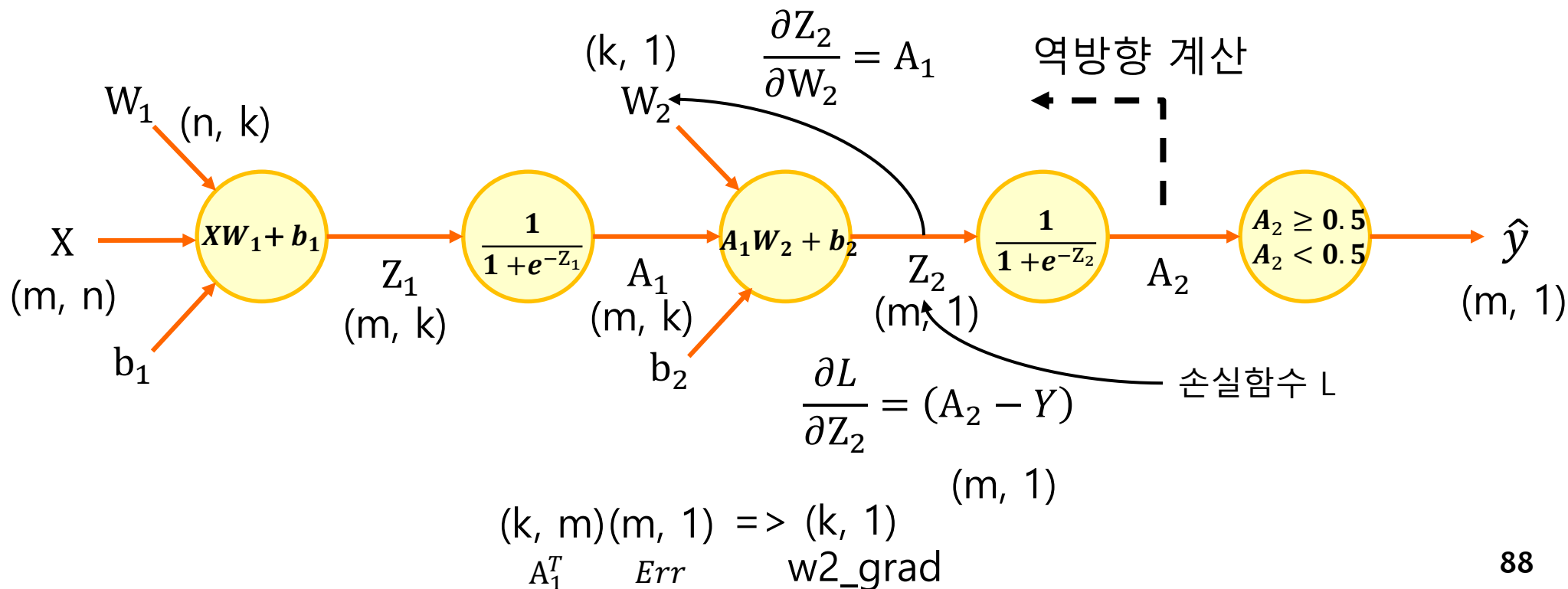
다층 신경망 backprop 이해

가중치 W_2 대하여 손실 함수를 미분한다.

$$\frac{\partial L}{\partial W_2} = \frac{\partial L}{\partial Z_2} \frac{\partial Z_2}{\partial W_2}$$

$$\frac{\partial L}{\partial W_2} = (A_1^T \cdot Err)$$

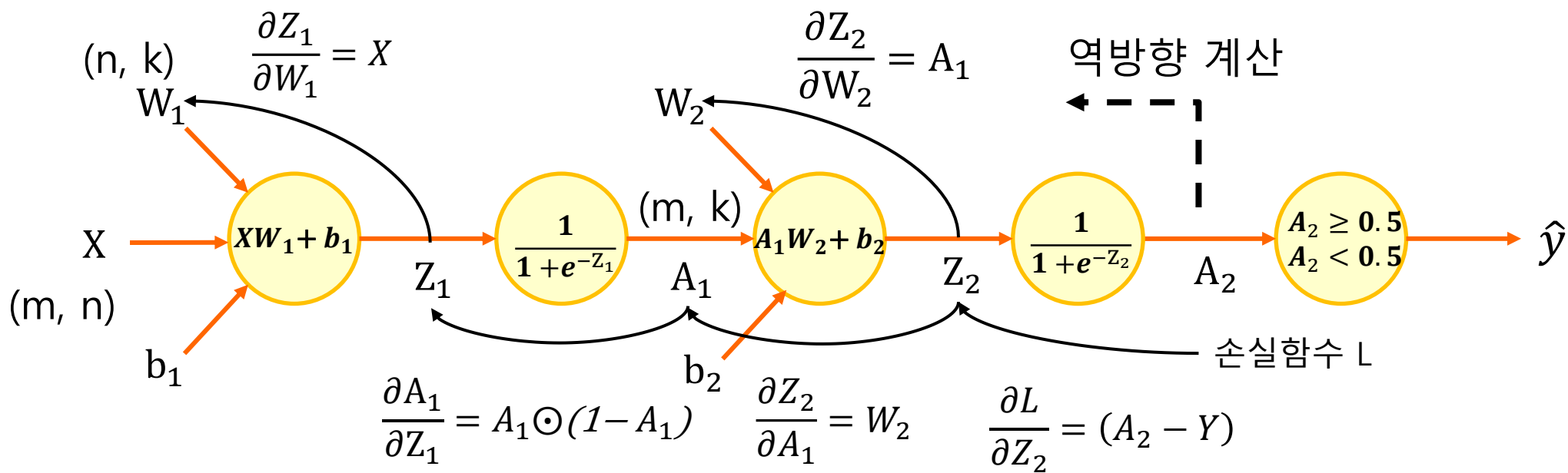
`np.dot(self.A1.T, Err)`



가중치 W_1 대하여 손실 함수를 미분한다. \odot : 멤버 곱 $a \odot (1-a)$

$$\frac{\partial L}{\partial W_1} = \frac{\partial L}{\partial Z_2} \frac{\partial Z_2}{\partial A_1} \frac{\partial A_1}{\partial Z_1} \frac{\partial Z_1}{\partial W_1} = X^T \left((A_2 - Y) W_2^T \odot A_1 \odot (1 - A_1) \right)$$

(m, 1)(1, k) => (m, k)
(n, m)(m, k) => (n, k)



미니배치 경사 하강법

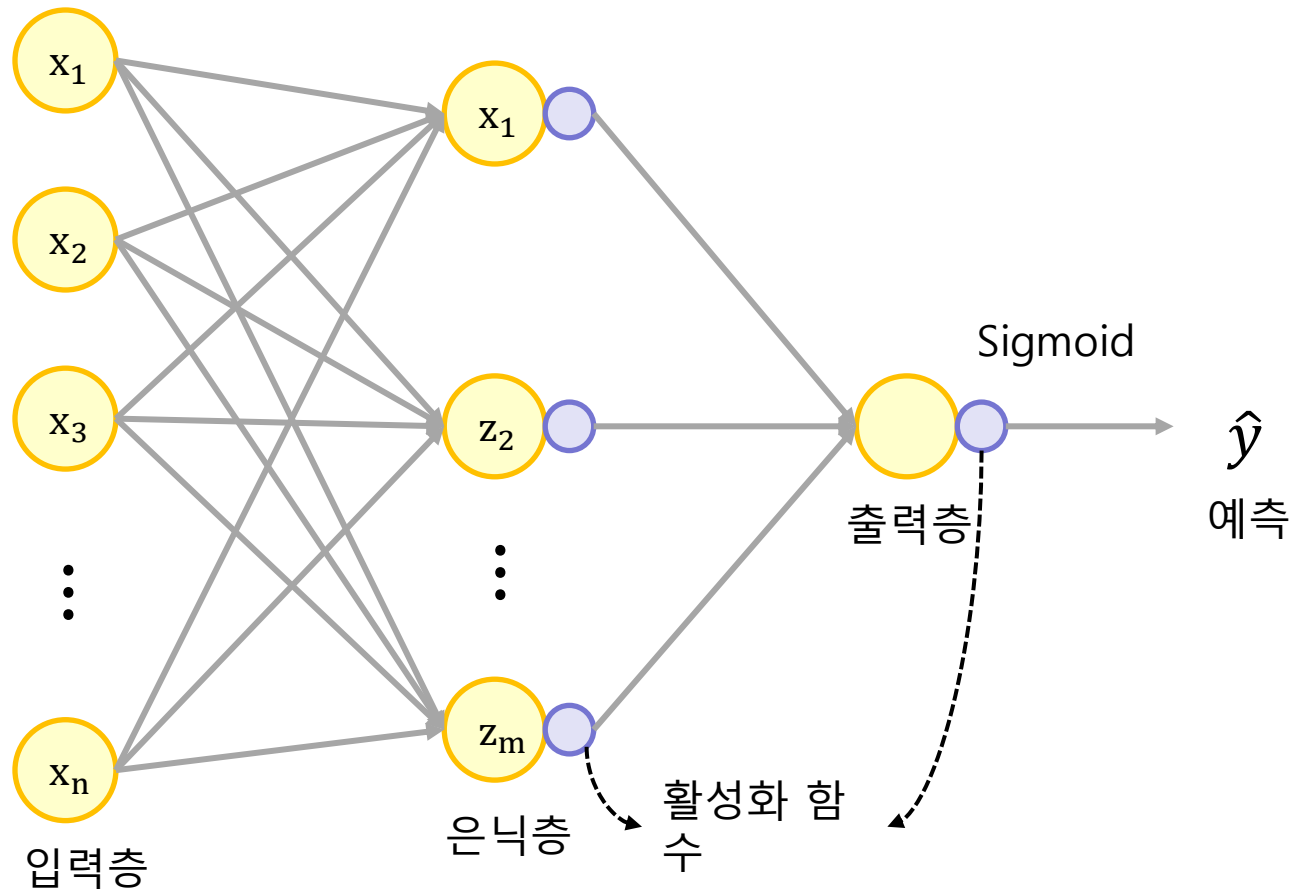
미니배치 경사 하강법 이란?

배치 경사 하강법과 비슷하지만 에포크마다 전체 데이터를 사용하는 것이 아니라 조금씩 나누어 계산을 수행하고 그레디언트를 구하여 가중치를 업데이트 한다.

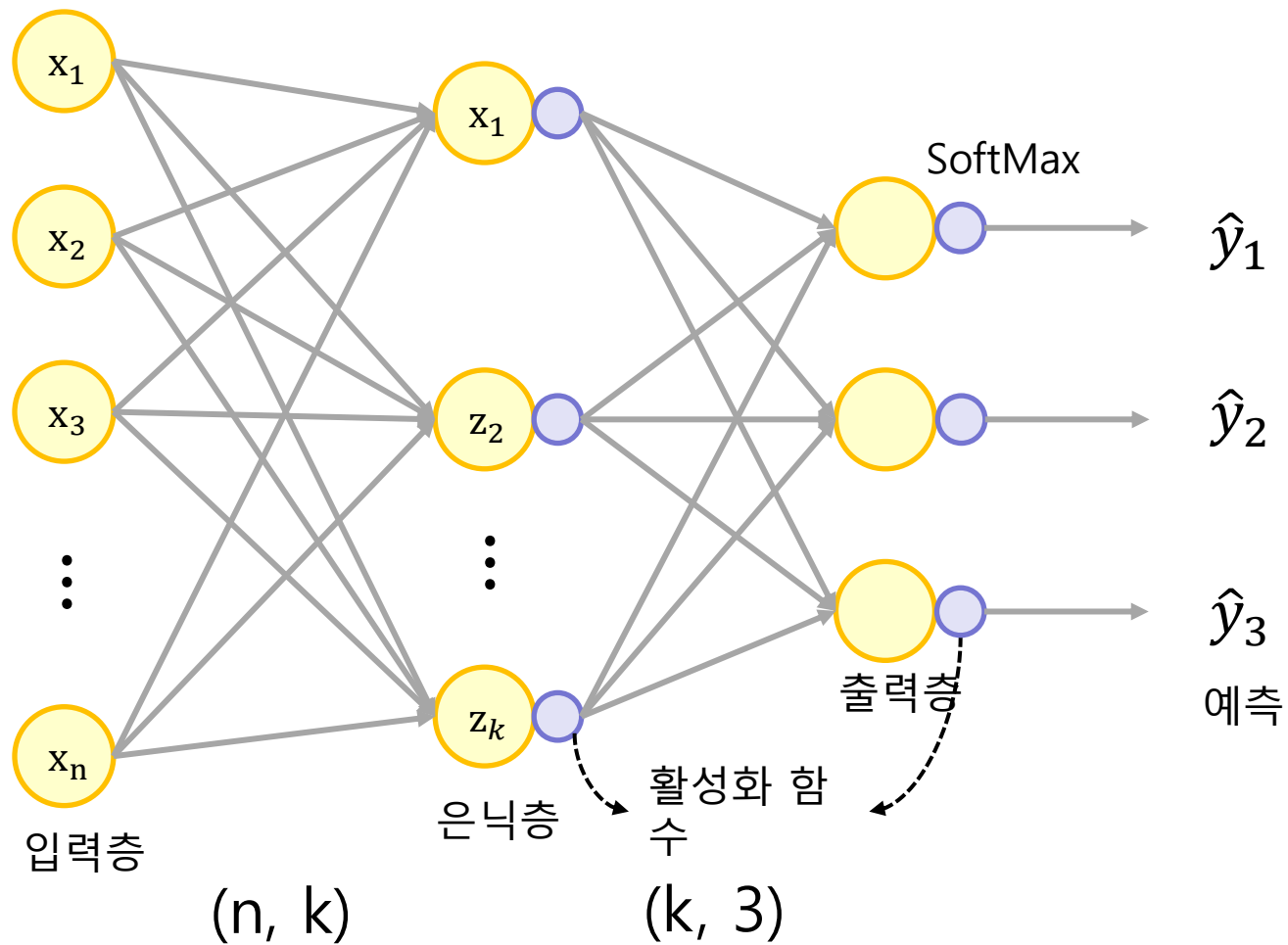
가중치 업데이트 방법

- 작게 나눈 미니 배치만큼 가중치를 업데이트 한다.
- 미니 배치의 크기는 보통 16,32,64등 2의 배수를 사용한다.
- 미니 배치의 크기가 1이면 확률적 경사 하강법이 된다.
- 미니 배치의 크기가 작으면 확률적 경사 하강법 처럼 작동하고 크면 배치 경사 하강법 처럼 작동한다.
- 미니 배치의 크기도 하이퍼파라미터이고 튜닝의 대상이다.

다중 분류 신경망



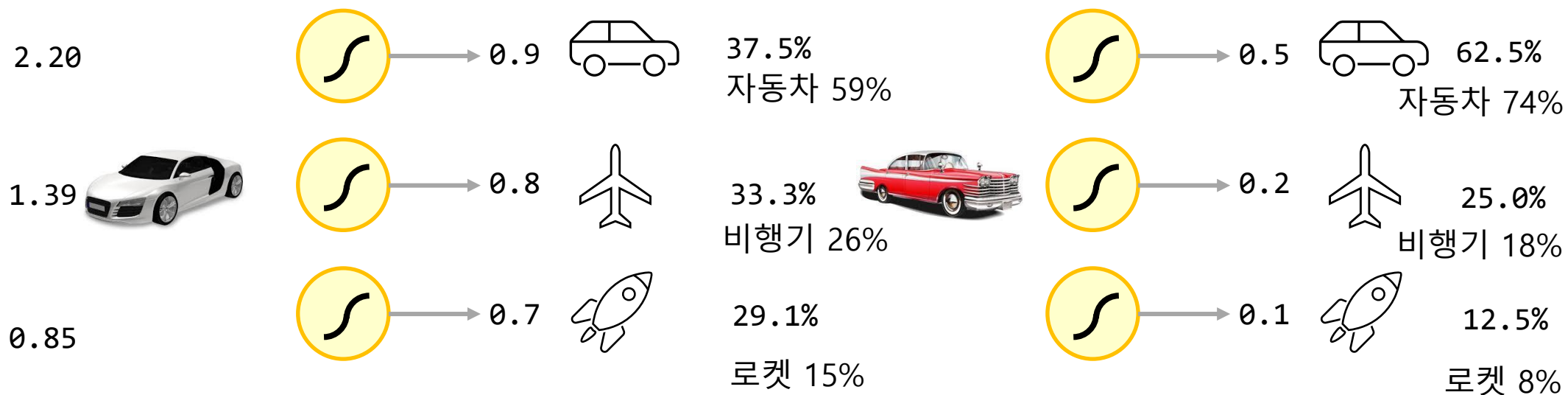
다중 분류 신경망



다중 분류 신경망

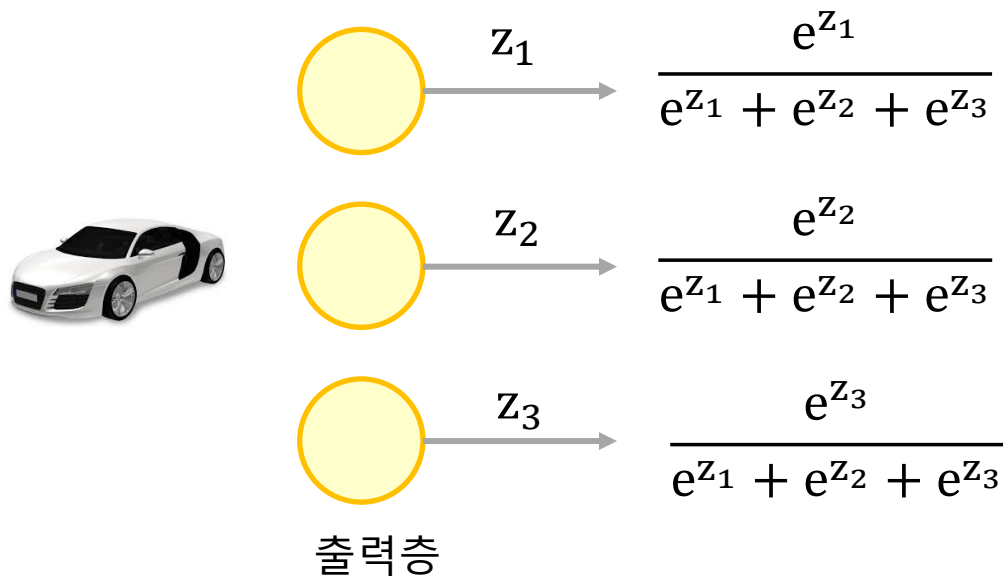
활성화 출력의 합이 1이 아니면 비교하기 어렵다.

소프트맥스 함수를 적용해 출력 강도를 정규화 한다.



소프트맥스 함수를 적용해 출력 강도를 정규화 한다.

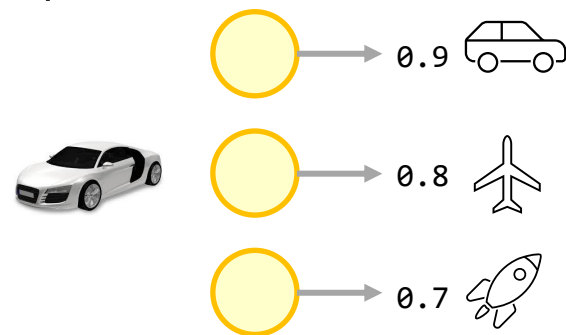
$$\frac{e^{z_i}}{e^{z_1} + e^{z_2} + e^{z_3}}$$



소프트맥스 함수를 적용해 출력 강도를 정규화 한다.

시그 모이드 함수를 z 에 대해 정리하면

$$\hat{y} = \frac{1}{1 + e^{-z}} \quad \Rightarrow \quad z = -\log\left(\frac{1}{\hat{y}} - 1\right)$$



$$z_1 = -\log\left(\frac{1}{0.9} - 1\right) = 2.20$$

$$z_2 = -\log\left(\frac{1}{0.8} - 1\right) = 1.39$$

$$z_3 = -\log\left(\frac{1}{0.7} - 1\right) = 0.85$$

$$\hat{y} = \frac{e^{2.20}}{e^{2.20} + e^{1.39} + e^{0.85}} = 0.59$$

$$\hat{y} = \frac{e^{1.39}}{e^{2.20} + e^{1.39} + e^{0.85}} = 0.26$$

$$\hat{y} = \frac{e^{0.85}}{e^{2.20} + e^{1.39} + e^{0.85}} = 0.15$$

자동차 59%

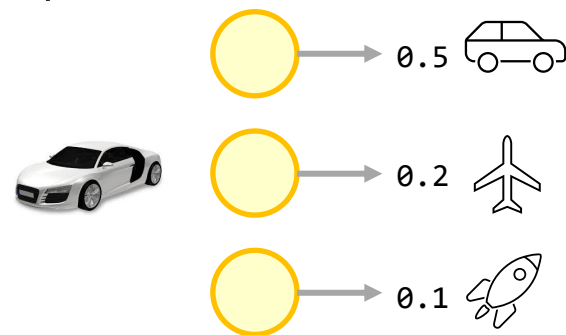
비행기 26%

로켓 15%

소프트맥스 함수를 적용해 출력 강도를 정규화 한다.

시그 모이드 함수를 z 에 대해 정리하면

$$\hat{y} = \frac{1}{1 + e^{-z}} \quad \Rightarrow \quad z = -\log\left(\frac{1}{\hat{y}} - 1\right)$$



$$z_1 = -\log\left(\frac{1}{0.5} - 1\right) = 0.00$$

$$z_2 = -\log\left(\frac{1}{0.2} - 1\right) = -1.39$$

$$z_3 = -\log\left(\frac{1}{0.1} - 1\right) = -2.20$$

$$\hat{y} = \frac{e^{0.00}}{e^{0.00} + e^{-1.39} + e^{-2.20}} = 0.74$$

$$\hat{y} = \frac{e^{-1.39}}{e^{0.00} + e^{-1.39} + e^{-2.20}} = 0.18$$

$$\hat{y} = \frac{e^{-2.20}}{e^{0.00} + e^{-1.39} + e^{-2.20}} = 0.08$$

자동차 74%

비행기 18%

로켓 8%

크로스 엔트로피 손실 함수

$$L = - \sum_{c=1}^c y_c \log(a_c) = -(y_1 \log(a_1) + y_2 \log(a_2) + \cdots + y_c \log(a_c))$$

크로스 엔트로피 손실 함수 미분

$$\frac{\partial L}{\partial z_1} = \frac{\partial L}{\partial a_1} \frac{\partial a_1}{\partial z_1} + \frac{\partial L}{\partial a_2} \frac{\partial a_2}{\partial z_1} + \frac{\partial L}{\partial a_3} \frac{\partial a_3}{\partial z_1}$$

크로스 엔트로피 손실 함수 미분

$$\frac{\partial L}{\partial z_1} = \frac{\partial L}{\partial a_1} \frac{\partial a_1}{\partial z_1} + \frac{\partial L}{\partial a_2} \frac{\partial a_2}{\partial z_1} + \frac{\partial L}{\partial a_3} \frac{\partial a_3}{\partial z_1}$$

$$\frac{\partial L}{\partial a_1} = -\frac{y_1}{a_1} \quad \frac{\partial L}{\partial a_2} = -\frac{y_2}{a_2} \quad \frac{\partial L}{\partial a_3} = -\frac{y_3}{a_3}$$

$$\frac{\partial L}{\partial z_1} = \left(-\frac{y_1}{a_1}\right) \frac{\partial a_1}{\partial z_1} + \left(-\frac{y_2}{a_2}\right) \frac{\partial a_2}{\partial z_1} + \left(-\frac{y_3}{a_3}\right) \frac{\partial a_3}{\partial z_1}$$

크로스 엔트로피 손실 함수 미분

$$\frac{\partial L}{\partial z_1} = \left(-\frac{y_1}{a_1}\right) \frac{\partial a_1}{\partial z_1} + \left(-\frac{y_2}{a_2}\right) \frac{\partial a_2}{\partial z_1} + \left(-\frac{y_3}{a_3}\right) \frac{\partial a_3}{\partial z_1}$$

$$\frac{\partial a_1}{\partial z_1} = a_1(1-a_1) \quad \frac{\partial a_2}{\partial z_1} = -a_2 a_1 \quad \frac{\partial a_3}{\partial z_1} = -a_3 a_1$$

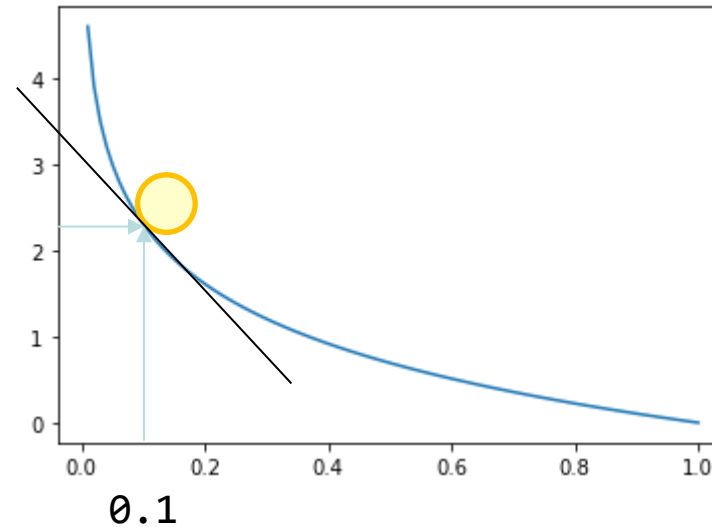
$$\begin{aligned} \frac{\partial L}{\partial z_1} &= \left(-\frac{y_1}{a_1}\right) a_1(1-a_1) + \left(-\frac{y_2}{a_2}\right) (-a_2 a_1) + \left(-\frac{y_3}{a_3}\right) (-a_3 a_1) \\ &= (a_1 - y_1) \end{aligned}$$

$$\frac{\partial L}{\partial z} = (\mathbf{a} - \mathbf{y})$$

	개	고양이	사슴
$a =$	$[0.7,$	$0.2,$	$0.1]$
$y =$	$[0,$	$0,$	$1]$

$$L = - \sum_{c=1}^3 y_c \log(a_c) = -(y_1 \log(a_1) + y_2 \log(a_2) + y_3 \log(a_3))$$

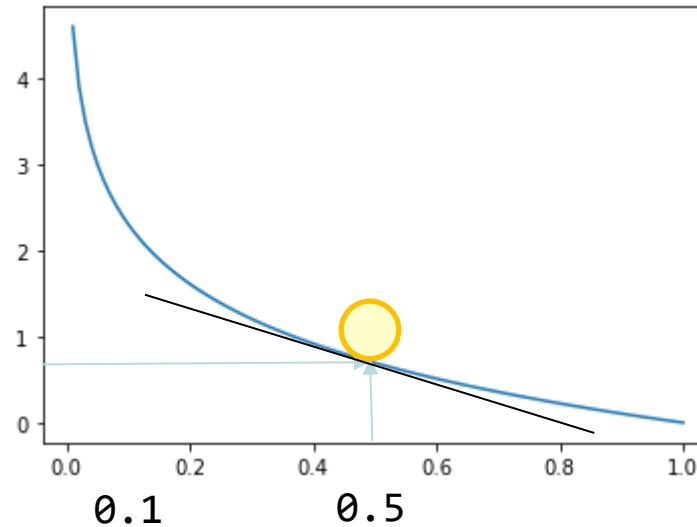
$$\begin{aligned} \text{err} &= (a - y) \\ \text{err} &= (0.1 - 1) = -0.9 \end{aligned}$$



	개	고양이	사슴
$a =$	$[0.1,$	$0.4,$	$0.5]$
$y =$	$[0,$	$0,$	$1]$

$$L = - \sum_{c=1}^3 y_c \log(a_c) = -(y_1 \log(a_1) + y_2 \log(a_2) + y_3 \log(a_3))$$

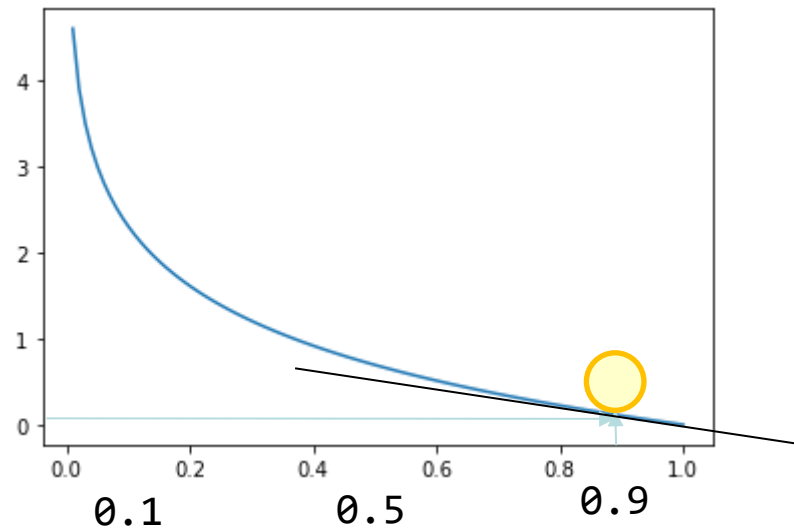
$$\begin{aligned} \text{err} &= (a - y) \\ \text{err} &= (0.5 - 1) = -0.5 \end{aligned}$$



	개	고양이	사슴
$a =$	$[0.0,$	$0.1,$	$0.9]$
$y =$	$[0,$	$0,$	$1]$

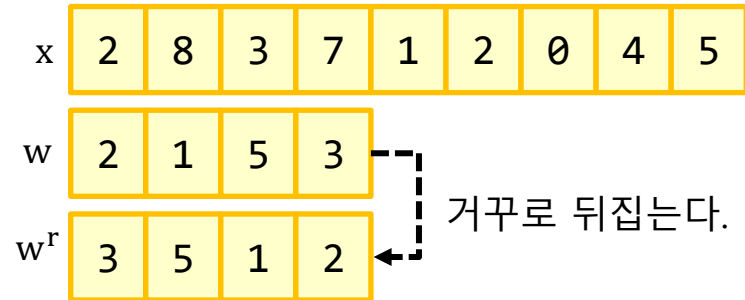
$$L = - \sum_{c=1}^3 y_c \log(a_c) = -(y_1 \log(a_1) + y_2 \log(a_2) + y_3 \log(a_3))$$

$$\begin{aligned} \text{err} &= (a - y) \\ \text{err} &= (0.9 - 1) = -0.1 \end{aligned}$$

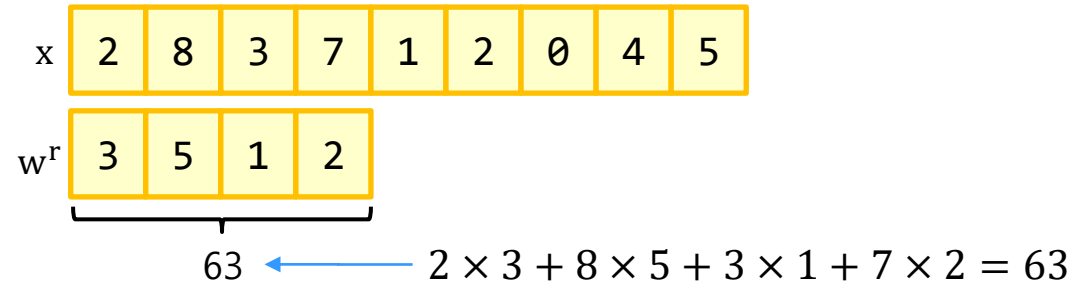


합성곱 연산

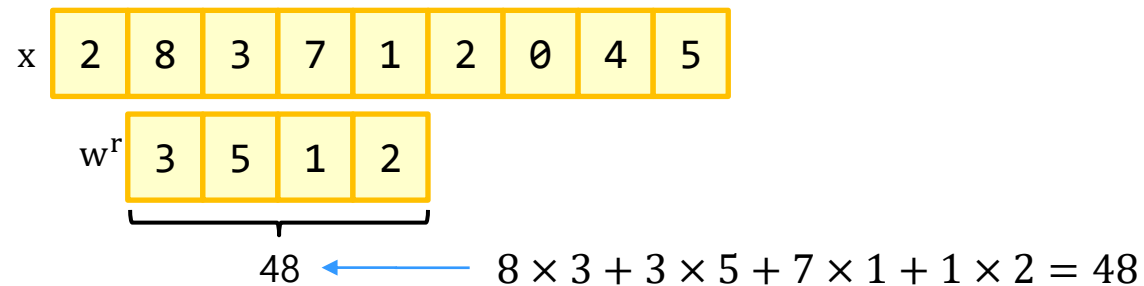
배열 하나 선택해 뒤집기



첫 번째 합성곱

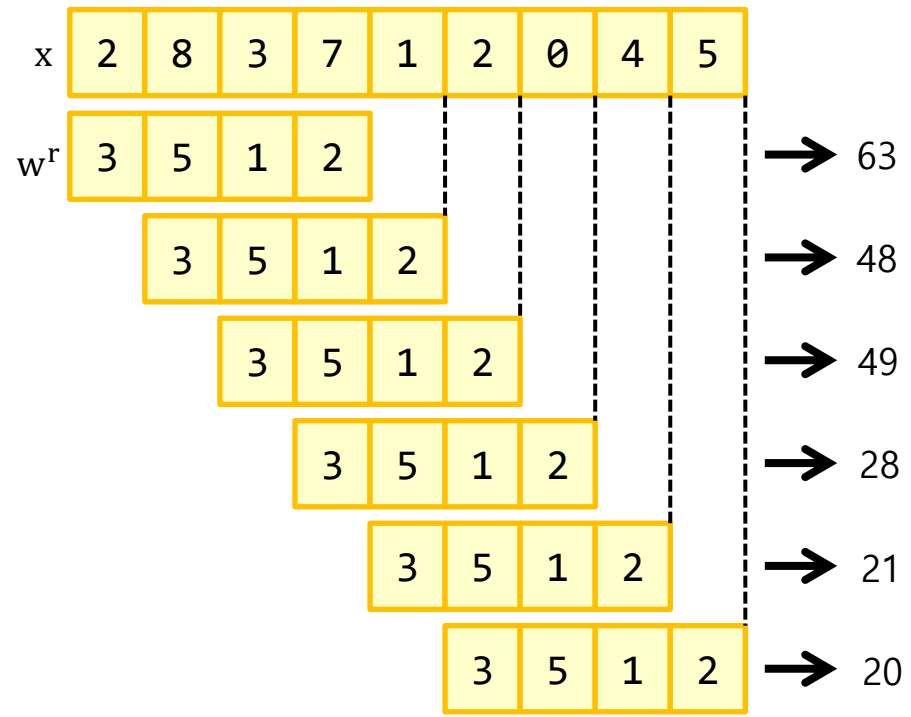


두 번째 합성곱



합성곱 연산

전체 합성곱



합성곱 구현

```
import numpy as np
x = np.array([2, 8, 3, 7, 1, 2, 0, 4, 5])
w = np.array([2, 1, 5, 3])
```

flip() 함수를 이용한 배열 뒤집기

```
w_r = np.flip(w)
print(w_r)
```

[3 5 1 2]

넘파이의 점 곱으로 합성곱 연산

```
for i in range(6):
    print(np.dot(x[i:i+4], w_r.reshape(-1,1)))
```

[63]
[48]
[49]
[28]
[21]
[20]

$$\begin{bmatrix} 2 & 8 & 3 & 7 \\ 8 & 3 & 7 & 1 \\ 3 & 7 & 1 & 2 \\ 7 & 1 & 2 & 0 \\ 1 & 2 & 0 & 4 \\ 2 & 0 & 4 & 5 \end{bmatrix} \cdot \begin{bmatrix} 3 \\ 5 \\ 1 \\ 2 \end{bmatrix} = \begin{bmatrix} 63 \\ 48 \\ 49 \\ 28 \\ 21 \\ 20 \end{bmatrix}$$

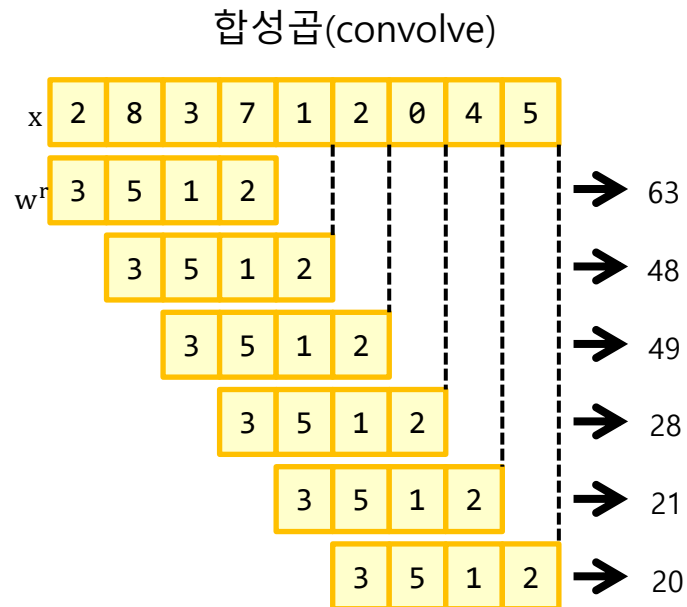
합성곱 연산

싸이파이로 합성곱 수행

```
from scipy.signal import convolve  
convolve(x, w, mode='valid')
```

```
array([63, 48, 49, 28, 21, 20])
```

합성곱 신경망은 진짜 합성곱을 사용하지 않는다.
합성곱 대신 교차상관을 사용한다.



합성곱 연산

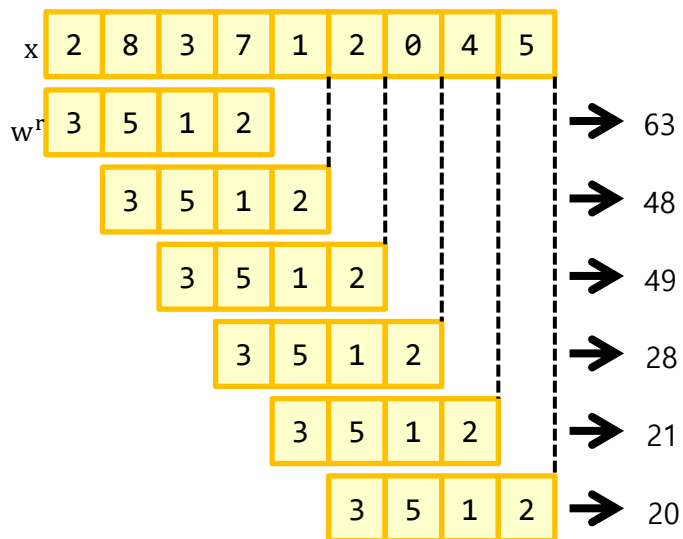
싸이파이로 교차상관 수행

```
from scipy.signal import correlate
correlate(x, w, mode='valid')
```

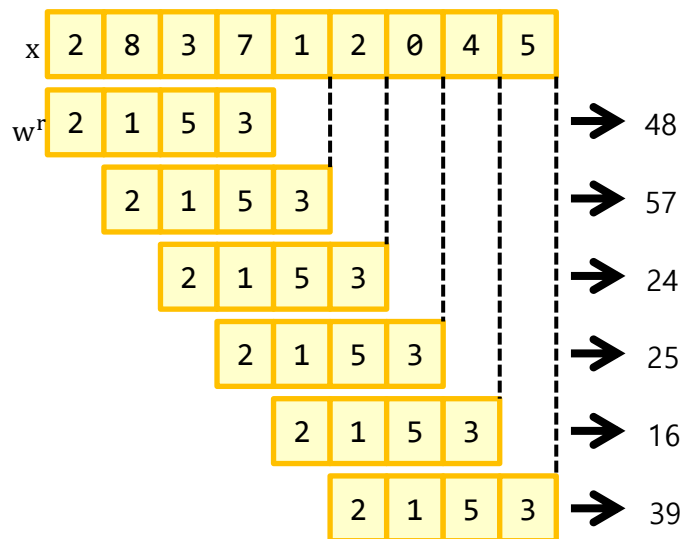
```
array([48, 57, 24, 25, 16, 39])
```

합성곱 신경망은 진짜 합성곱을 사용하지 않는다.
합성곱 대신 교차상관을 사용한다.

합성곱(convolve)



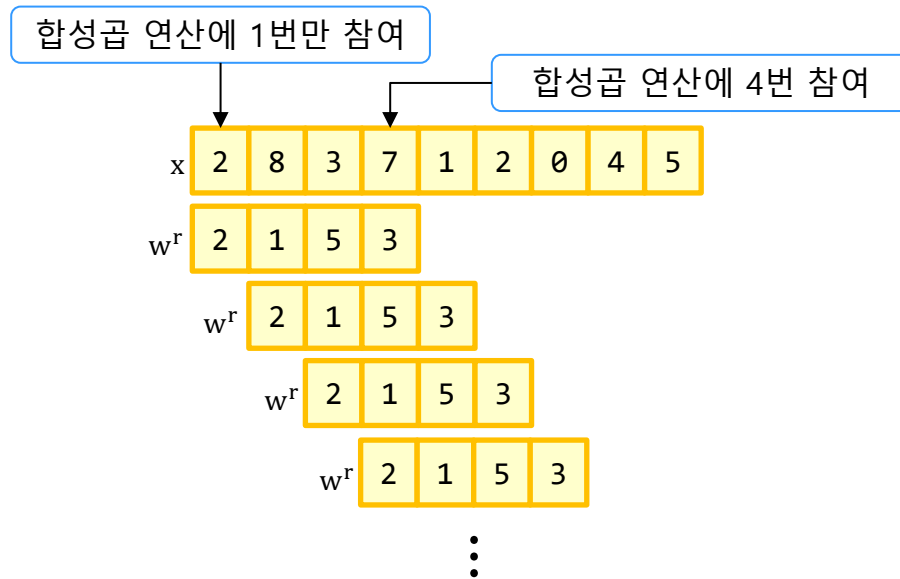
교차상관(correlate)



합성곱 연산

패딩과 스트라이드 이해

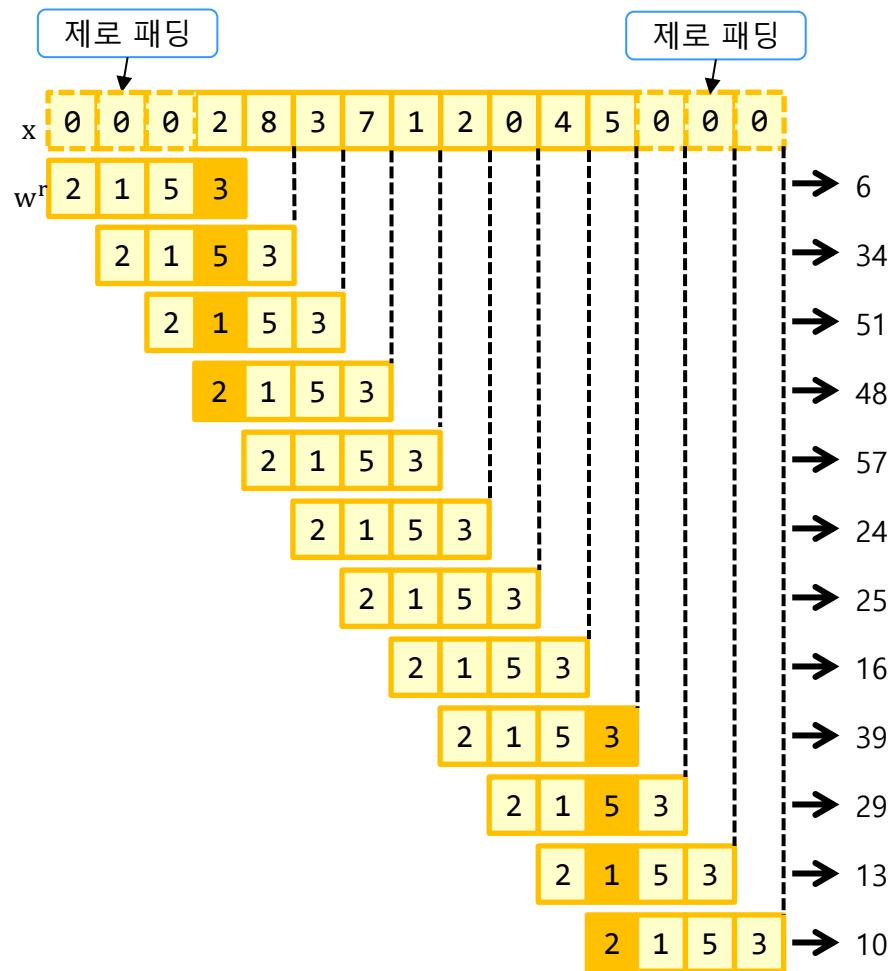
밸리드 패딩은 원본 배열의 원소가 합성곱 연산에 참여하는 정도가 다르다.



합성곱 연산

패딩과 스트라이드 이해

풀 패딩은 원본 배열의 원소의 연산 참여도를 동일하게 만든다.



```
correlate(x, w, mode='full')
```

```
array([ 6, 34, 51, 48, 57, 24, 25, 16, 39, 29, 13, 10])
```

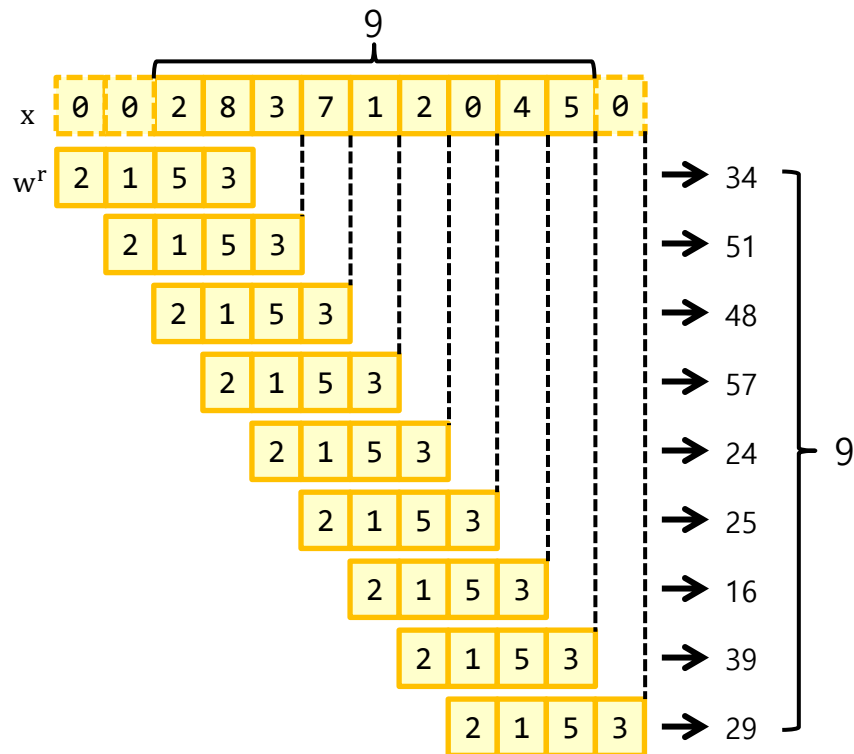
합성곱 연산

패딩과 스트라이드 이해

세임 패딩은 출력 배열의 길이를 원본 배열의 원소의 길이와 동일하게 만든다.

```
correlate(x, w, mode='same')
```

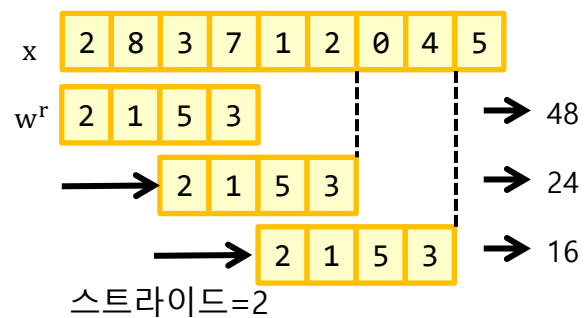
```
array([34, 51, 48, 57, 24, 25, 16, 39, 29])
```



합성곱 연산

패딩과 스트라이드 이해

스트라이드는 미끄러지는 간격을 조정한다.



```
correlate(x, w, mode='same')
```

```
array([34, 51, 48, 57, 24, 25, 16, 39, 29])
```

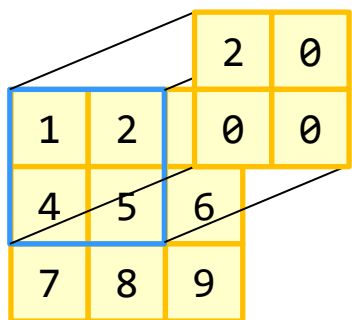
합성곱 연산

2차원 배열에서 합성곱 수행 (mode='valid')

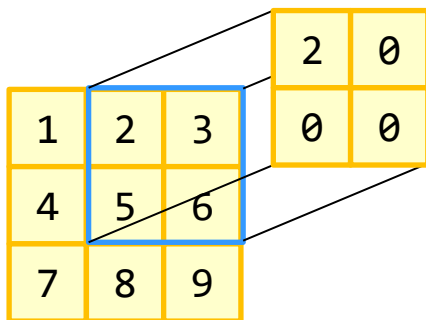
x			w	
1	2	3	2	0
4	5	6	0	0
7	8	9		

```
x = np.array([[1, 2, 3],
              [4, 5, 6],
              [7, 8, 9]])
w = np.array([[2, 0],
              [0, 0]])
from scipy.signal import correlate2d
correlate2d(x, w, mode='valid')
```

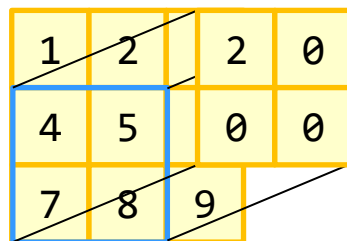
```
array([[ 2,  4],
       [ 8, 10]])
```



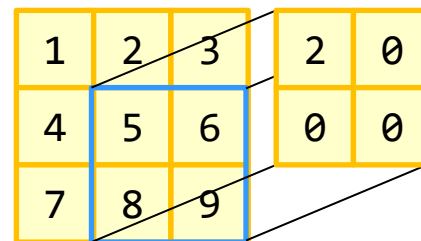
=> 2



=> 4



=> 8



=> 10

합성곱 연산

2차원 배열에서 same padding

1	2	3	0
4	5	6	0
7	8	9	0
0	0	0	0

1	2	3	0
4	5	6	0
7	8	9	0
0	0	0	0

1	2	3	0
4	5	6	0
7	8	9	0
0	0	0	0

1	2	3	0
4	5	6	0
7	8	9	0
0	0	0	0

1	2	3	0
4	5	6	0
7	8	9	0
0	0	0	0

1	2	3	0
4	5	6	0
7	8	9	0
0	0	0	0

1	2	3	0
4	5	6	0
7	8	9	0
0	0	0	0

1	2	3	0
4	5	6	0
7	8	9	0
0	0	0	0

1	2	3	0
4	5	6	0
7	8	9	0
0	0	0	0

1	2	3	0
4	5	6	0
7	8	9	0
0	0	0	0

1	2	3	0
4	5	6	0
7	8	9	0
0	0	0	0

1	2	3	0
4	5	6	0
7	8	9	0
0	0	0	0

1	2	3	0
4	5	6	0
7	8	9	0
0	0	0	0

1	2	3	0
4	5	6	0
7	8	9	0
0	0	0	0

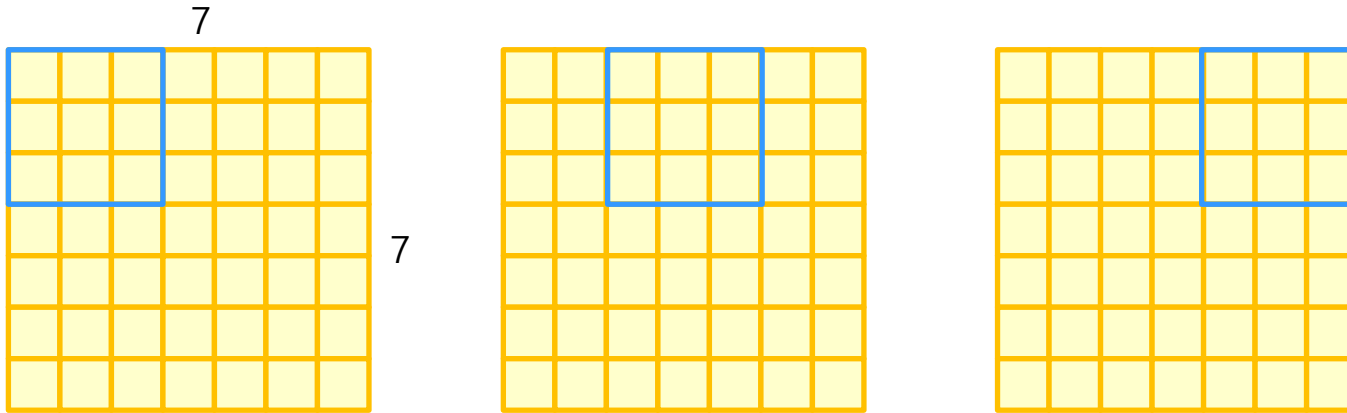
1	2	3	0
4	5	6	0
7	8	9	0
0	0	0	0

```
correlate2d(x, w, mode='same')
```

```
array([[ 2,  4,  6],  
       [ 8, 10, 12],  
       [14, 16, 18]])
```

합성곱 연산

2차원 배열에서 스트라이드 이해



7x7 input

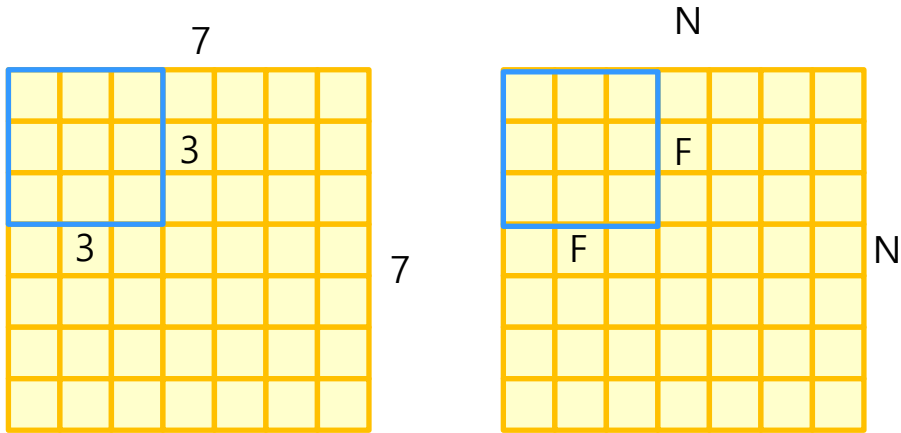
3x3 filter

stride = 2

=> 3x3 output

합성곱 연산

2차원 배열에서 스트라이드 이해



Output size :

$$(N - F) / \text{stride} + 1$$

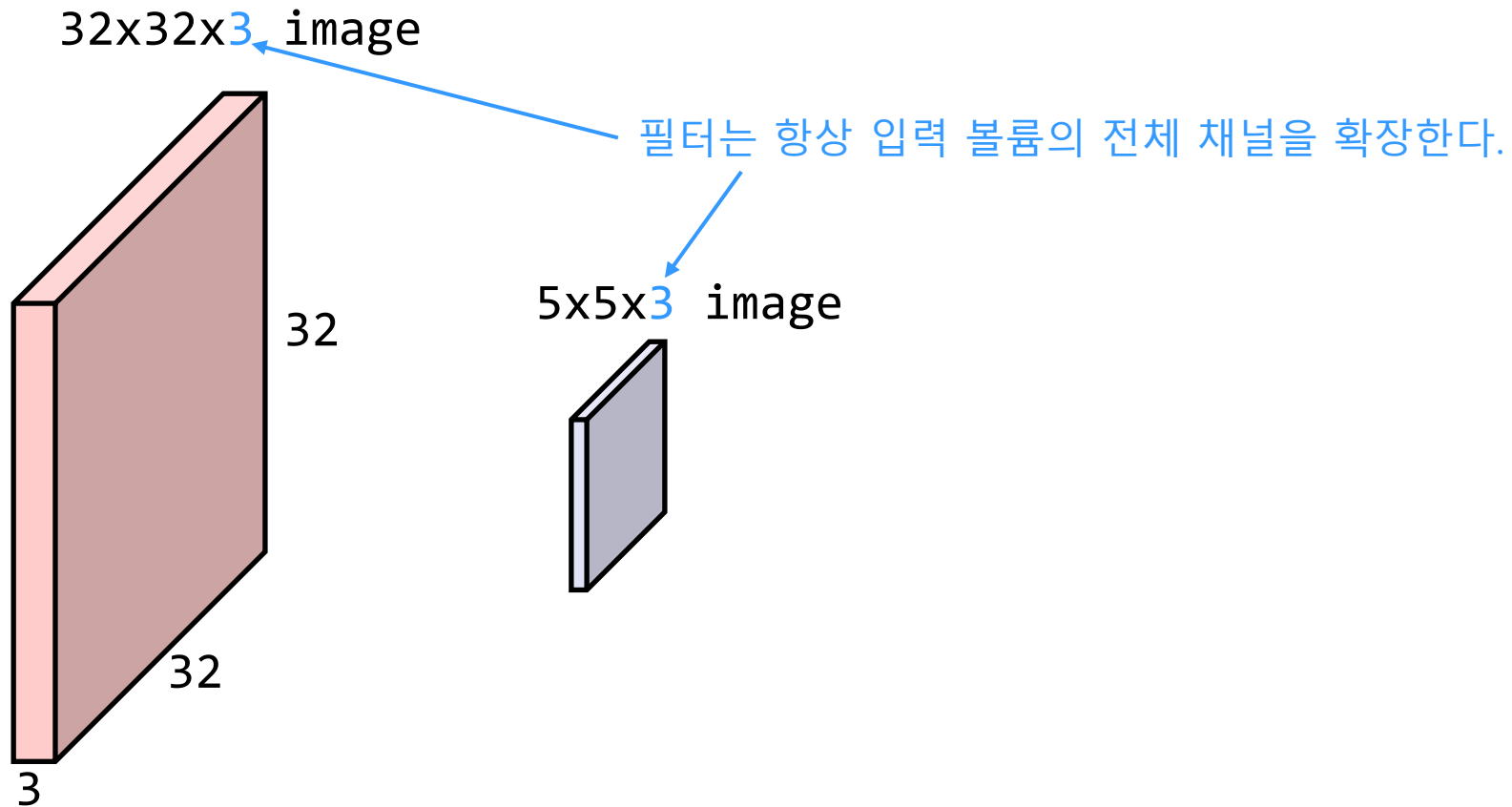
예) $N = 7, F = 3$

$$\text{stride } 1 \Rightarrow (7-3)/1+1 = 5$$

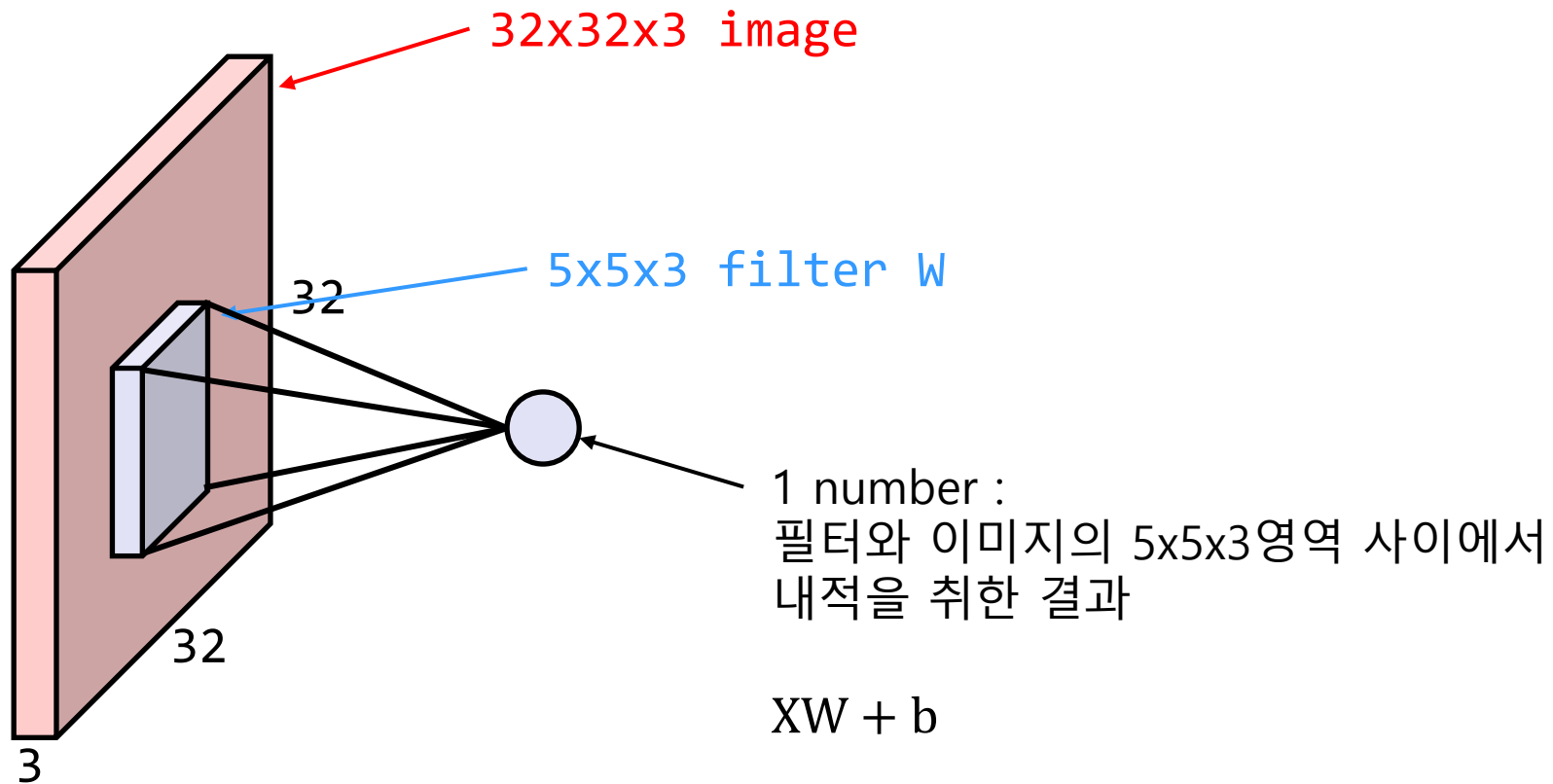
$$\text{stride } 2 \Rightarrow (7-3)/2+1 = 3$$

$$\text{stride } 3 \Rightarrow (7-3)/3+1 = 2.33$$

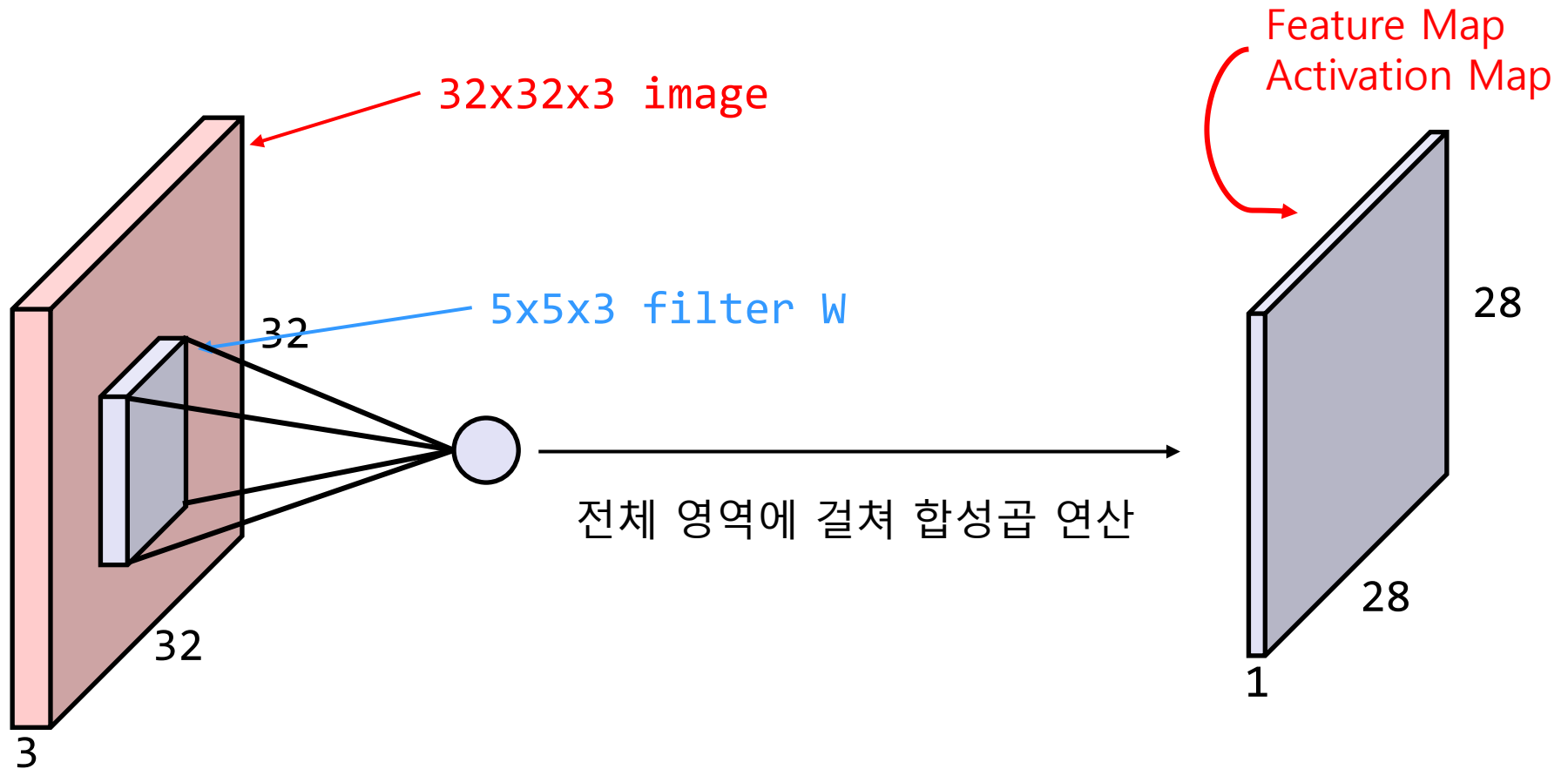
Convolution Layer



Convolution Layer

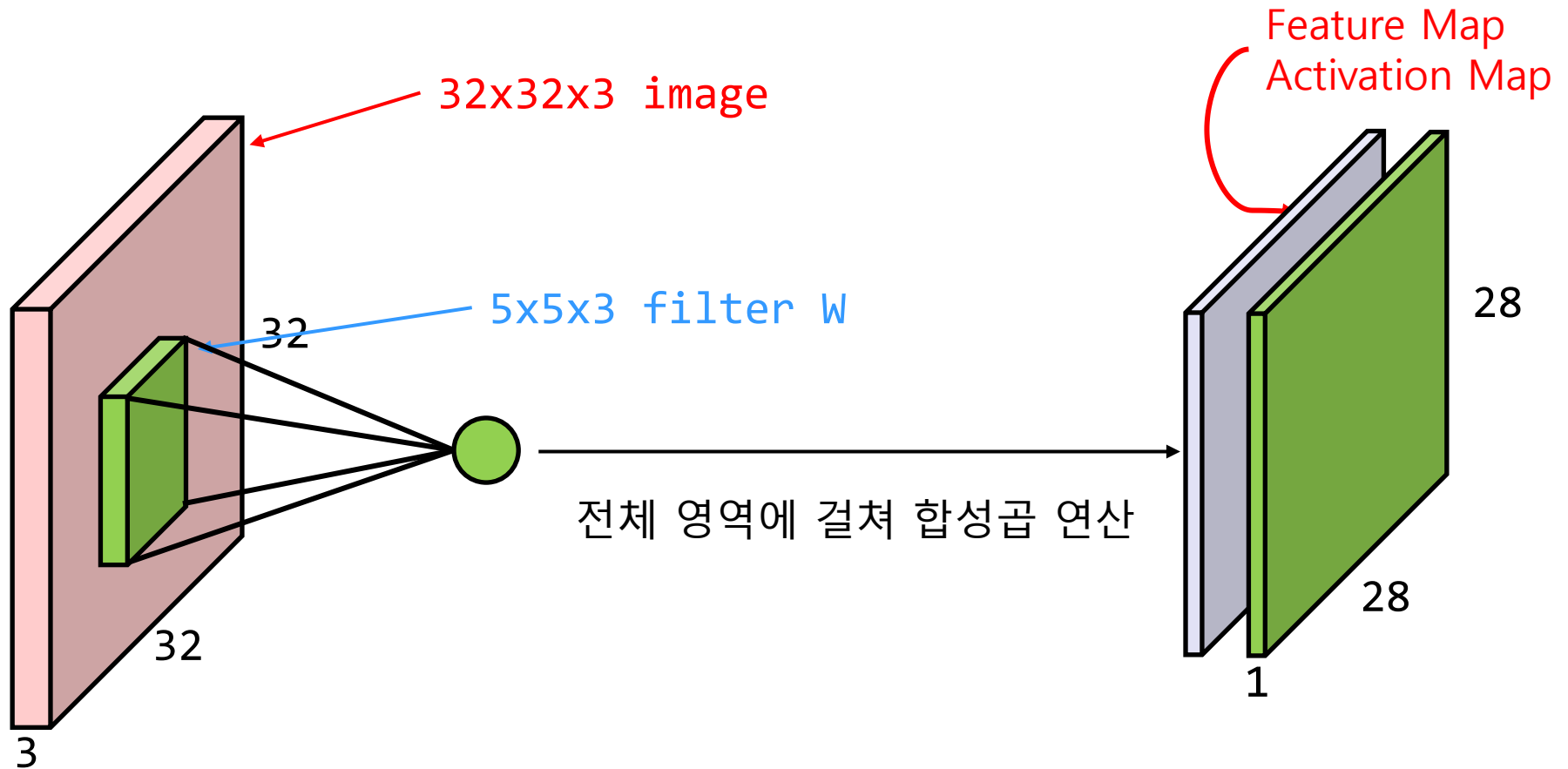


Convolution Layer



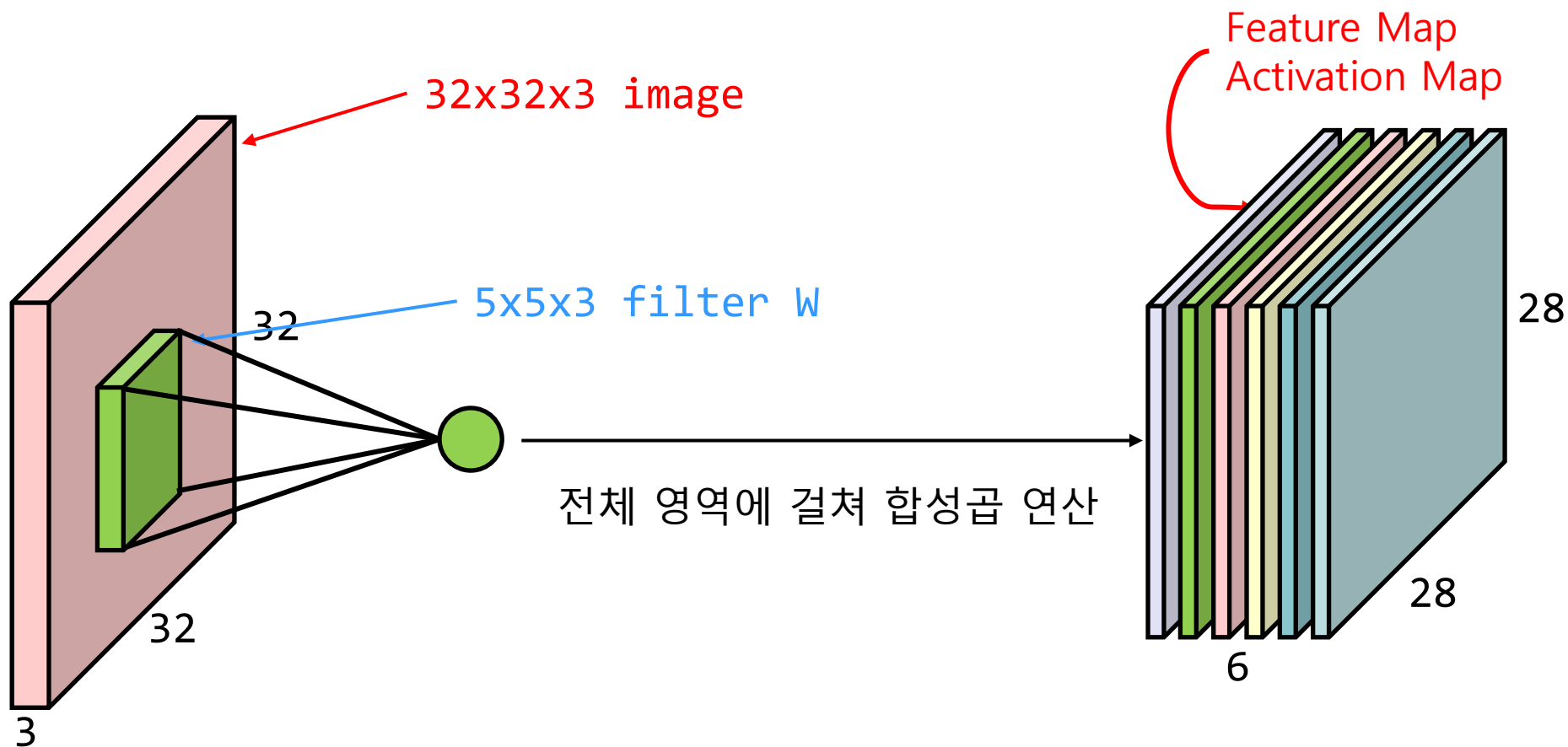
두번째 필터 동작

Convolution Layer



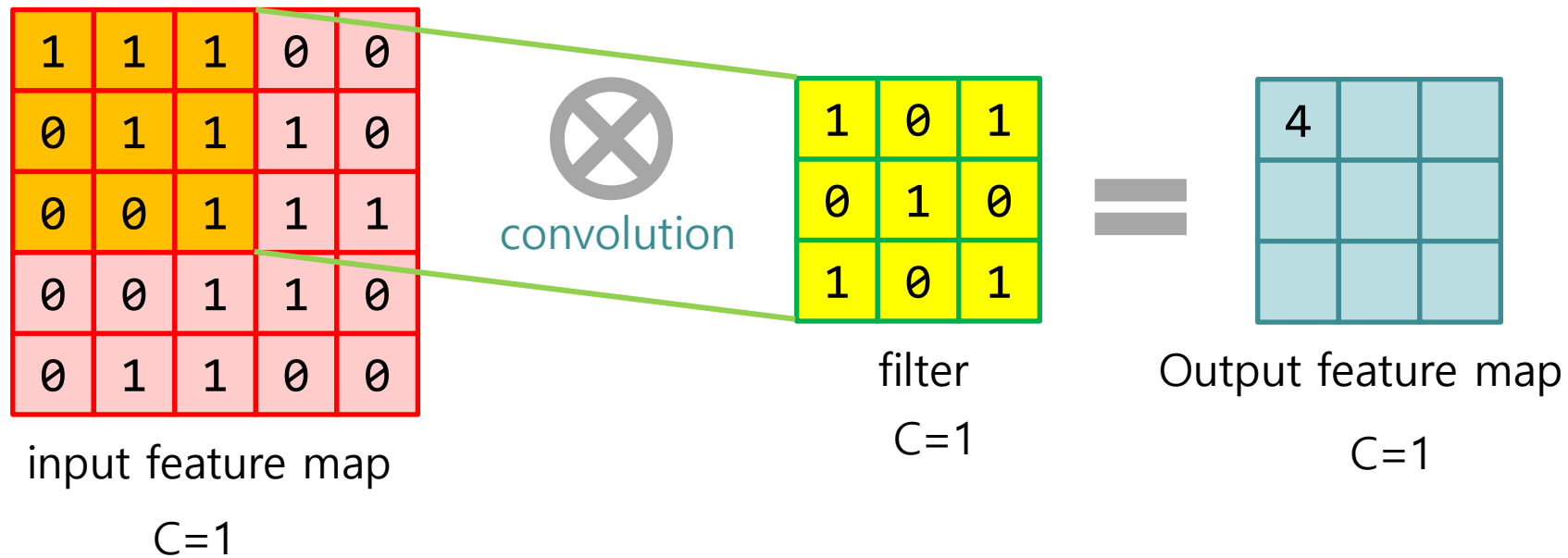
Convolution Layer

6개의 5x5필터가 있다면, 6개의 개별 feature map이 생성됨

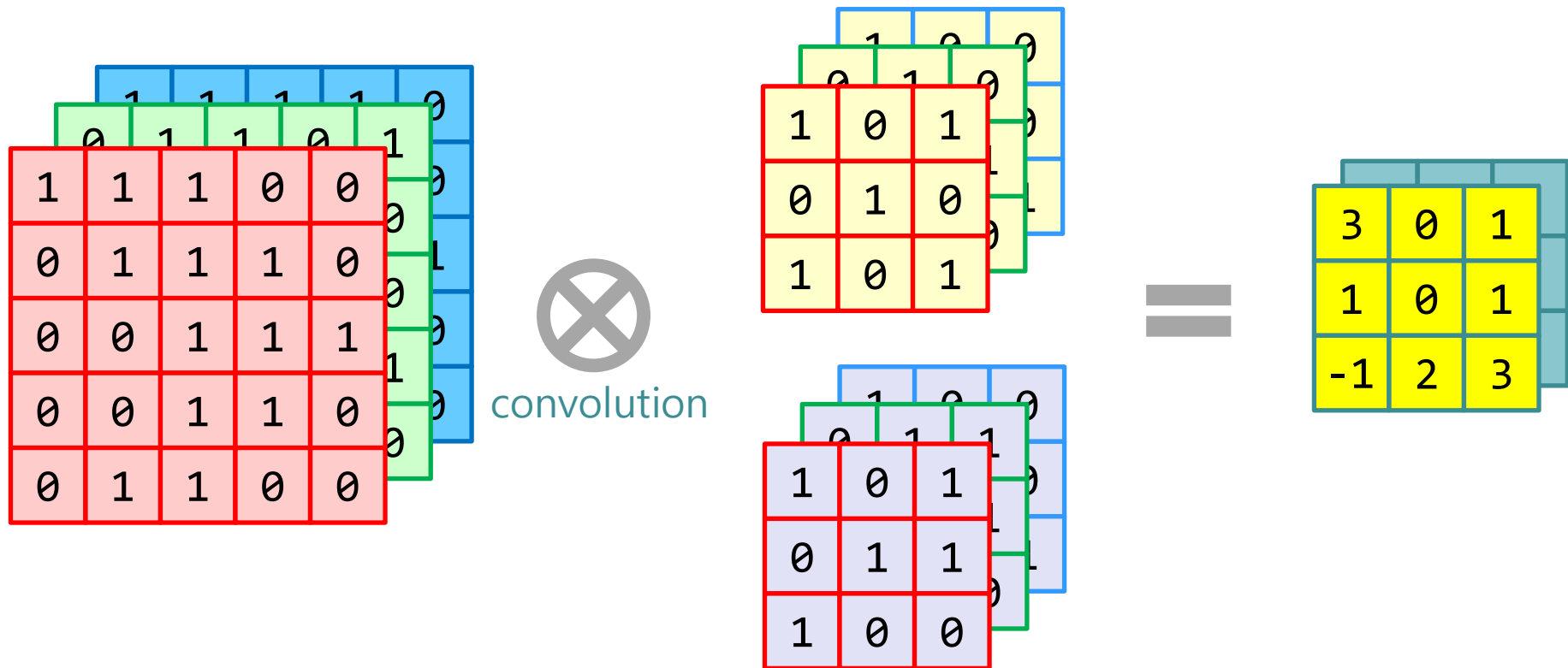


Convolution Layer - 계산

$$1 \times 1 + 1 \times 0 + 1 \times 1 + 0 \times 0 + 1 \times 1 + 1 \times 0 + 0 \times 1 + 0 \times 0 + 1 \times 1 = 4$$



Convolution Layer - 계산

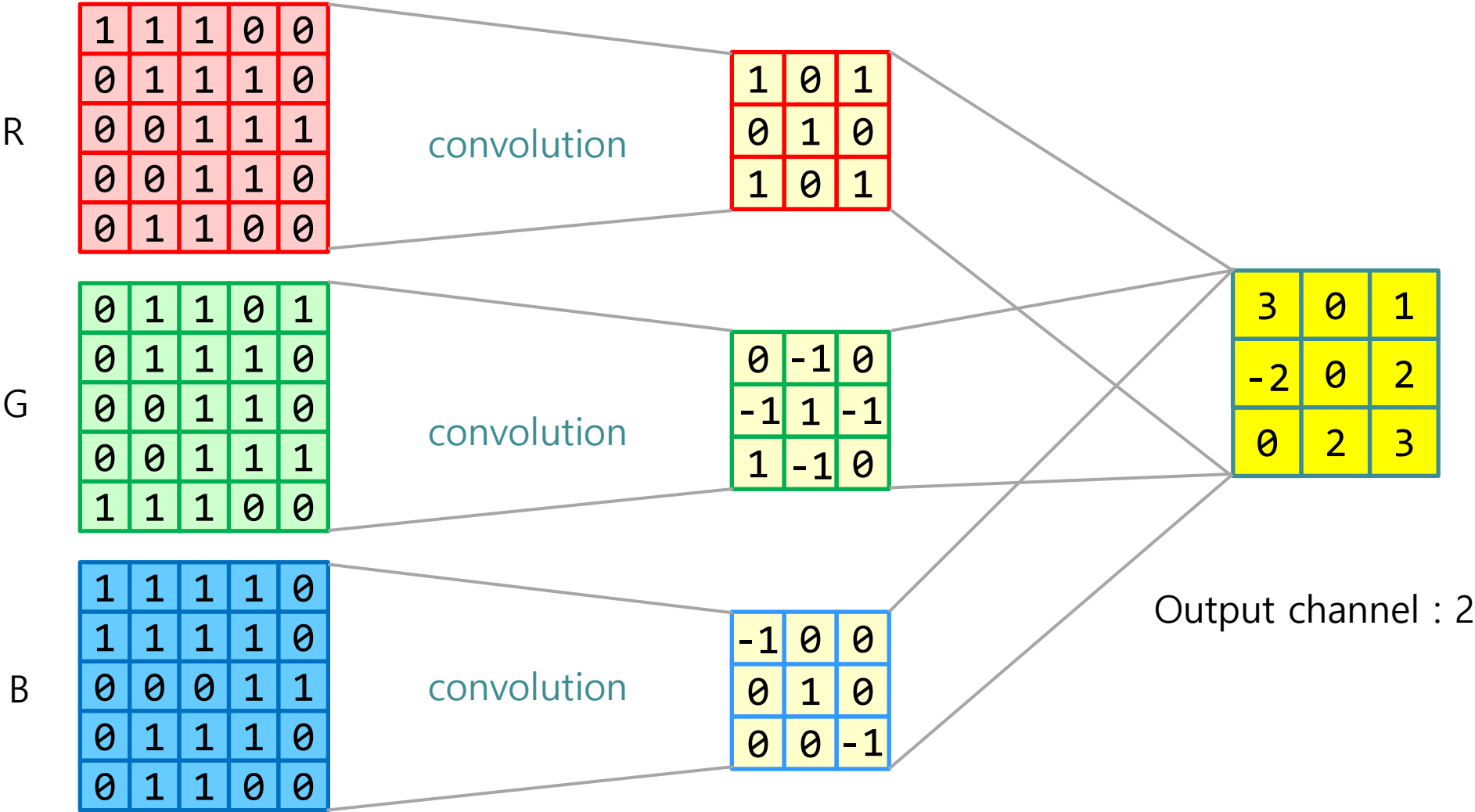


input channel = 3

of filters : 2

Output channel : 2

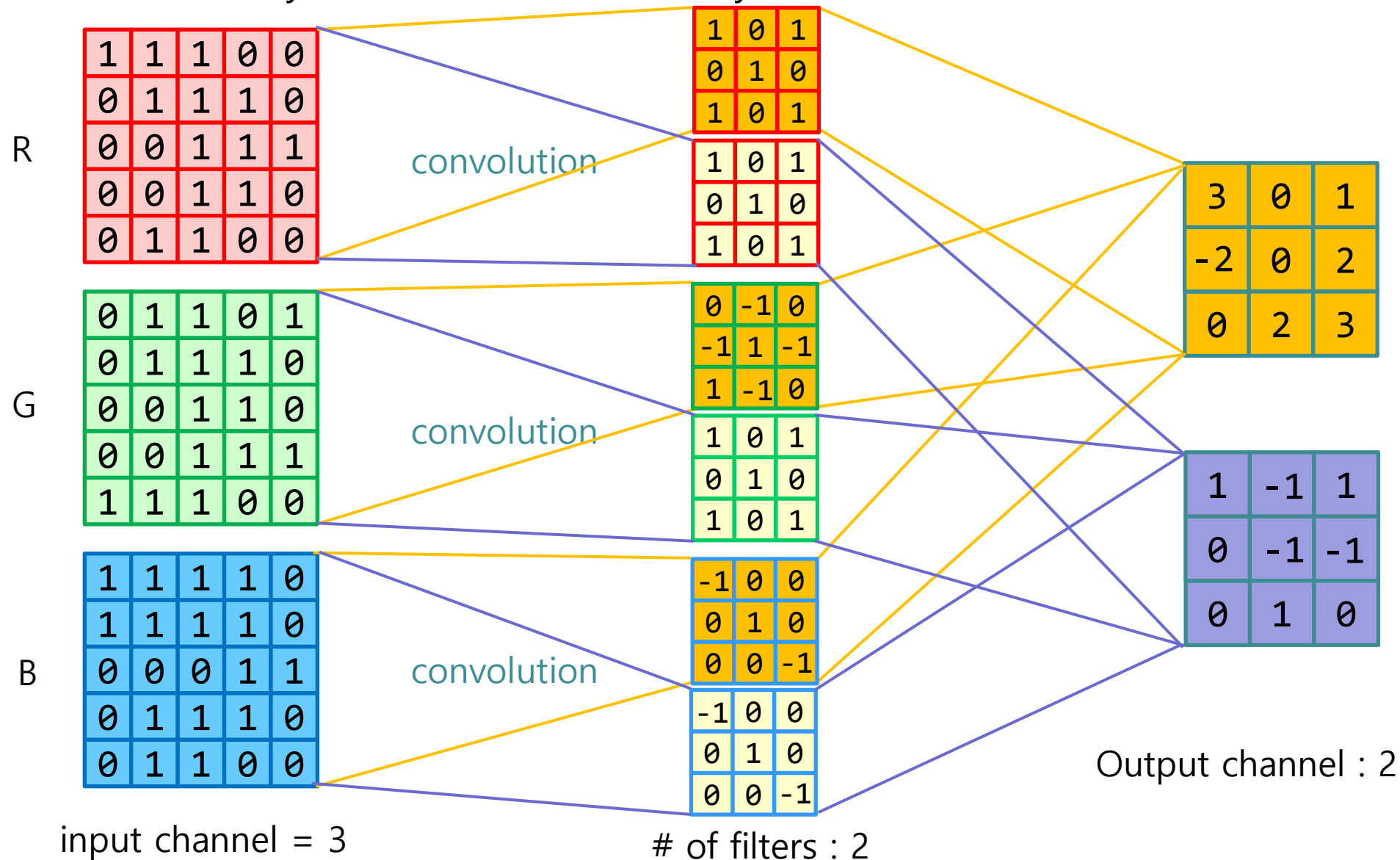
Convolution Layer - Multi Channel, Many Filters



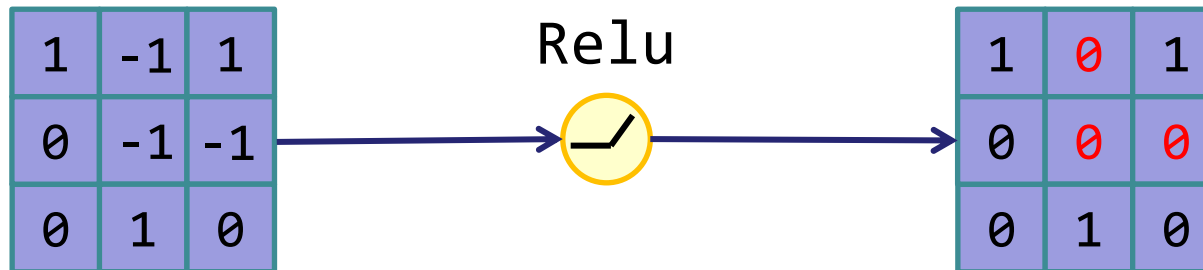
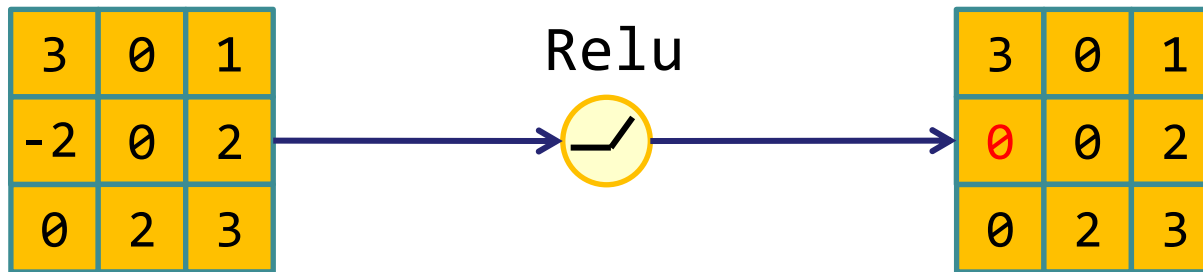
input channel = 3

of filters : 2

Convolution Layer - Multi Channel, Many Filters



Activation Function



tf.keras.layers.Conv2D

Arguments:

- **filters:** Integer, the dimensionality of the output space (i.e. the number of output filters in the convolution).
- **kernel_size:** An integer or tuple/list of 2 integers, specifying the height and width of the 2D convolution window. Can be a single integer to specify the same value for all spatial dimensions.
- **strides:** An integer or tuple/list of 2 integers, specifying the strides of the convolution along the height and width. Can be a single integer to specify the same value for all spatial dimensions. Specifying any stride value $\neq 1$ is incompatible with specifying any `dilation_rate` value $\neq 1$.
- **padding:** one of "valid" or "same" (case-insensitive).
- **data_format:** A string, one of `channels_last` (default) or `channels_first`. The ordering of the dimensions in the inputs. `channels_last` corresponds to inputs with shape (batch_size, height, width, channels) while `channels_first` corresponds to inputs with shape (batch_size, channels, height, width). It defaults to the `image_data_format` value found in your Keras config file at `~/.keras/keras.json`. If you never set it, then it will be "channels_last".

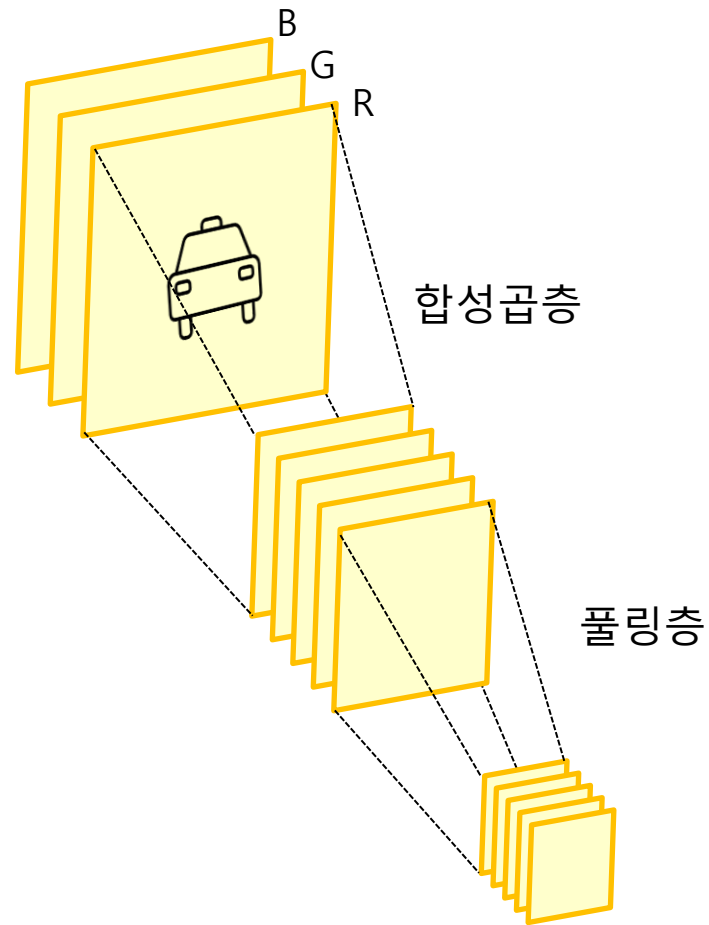
tf.keras.layers.Conv2D

- **activation:** Activation function to use. If you don't specify anything, no activation is applied (see [keras.activations](#)).
- **use_bias:** Boolean, whether the layer uses a bias vector.
- **kernel_initializer:** Initializer for the kernel weights matrix (see [keras.initializers](#)).
- **bias_initializer:** Initializer for the bias vector (see [keras.initializers](#)).
- **kernel_regularizer:** Regularizer function applied to the kernel weights matrix (see [keras.regularizers](#)).
- **bias_regularizer:** Regularizer function applied to the bias vector (see [keras.regularizers](#)).

풀링 연산

풀링 연산

합성곱층과 풀링층을 거치면서 변환되는 과정

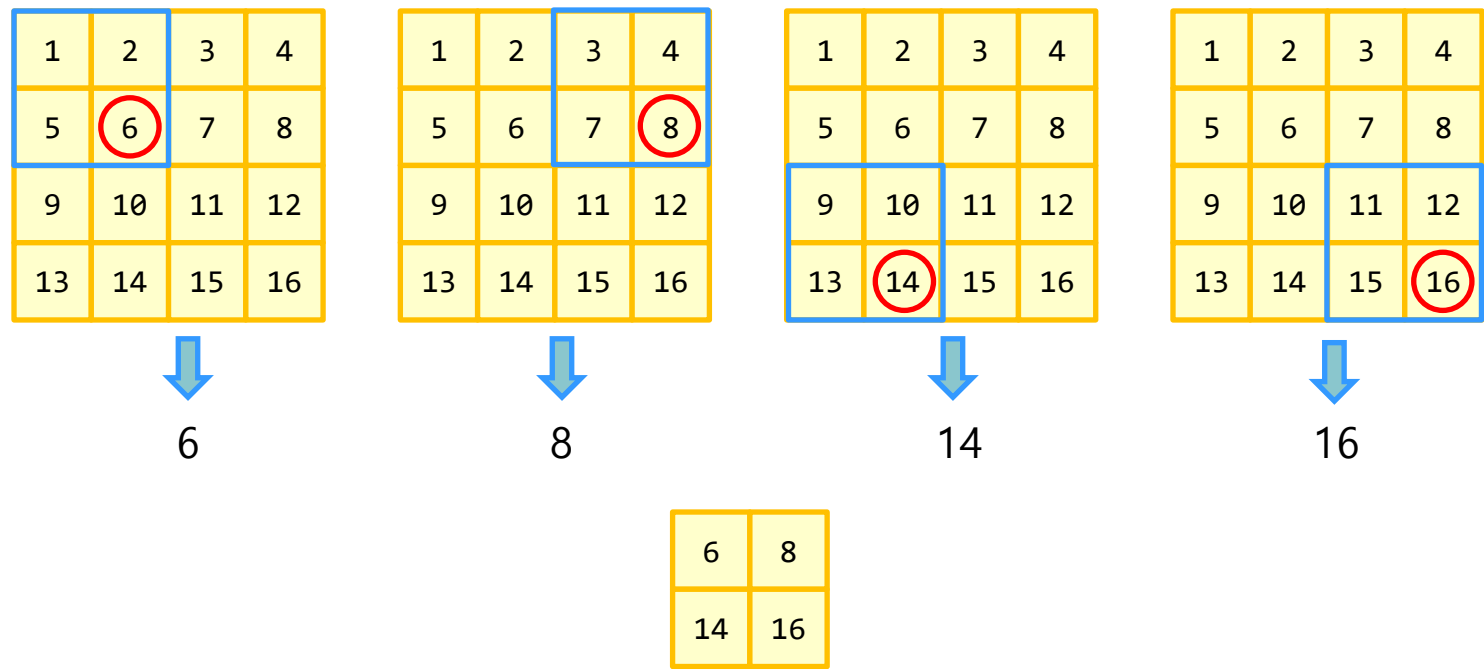


풀링 연산

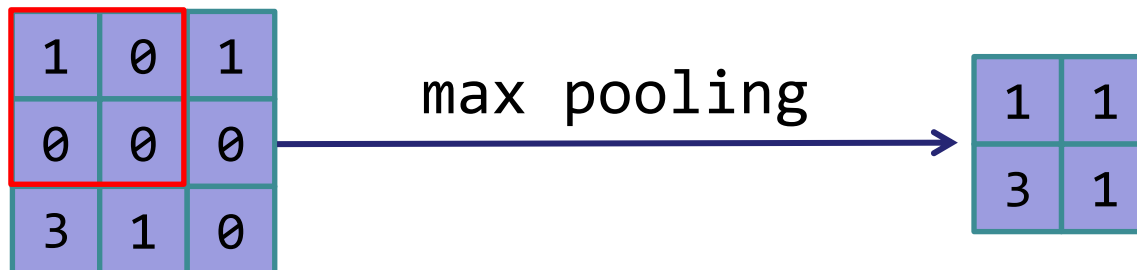
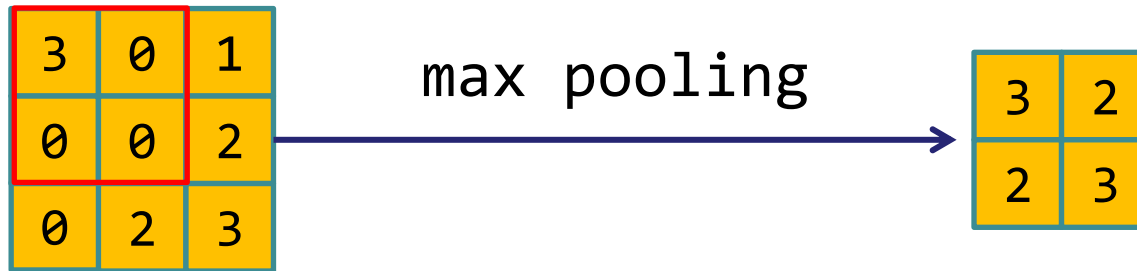
풀링 연산

풀링이란? 특성 맵을 스캔하며 최대값을 고르거나 평균값을 계산하는 것을 말함

최대 풀링



Pooling(max pooling, 2x2 filter, stride 1)



Pooling(average pooling, 2x2 filter, stride 2)

1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16



$$\frac{1 + 2 + 5 + 6}{4} = 3.5$$

1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16



$$\frac{3 + 4 + 7 + 8}{4} = 5.5$$

1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16



$$\frac{9 + 10 + 13 + 14}{4} = 11.5$$

1	2	3	4
5	6	7	8
9	10	11	12
13	14	15	16



$$\frac{11 + 12 + 15 + 16}{4} = 13.5$$

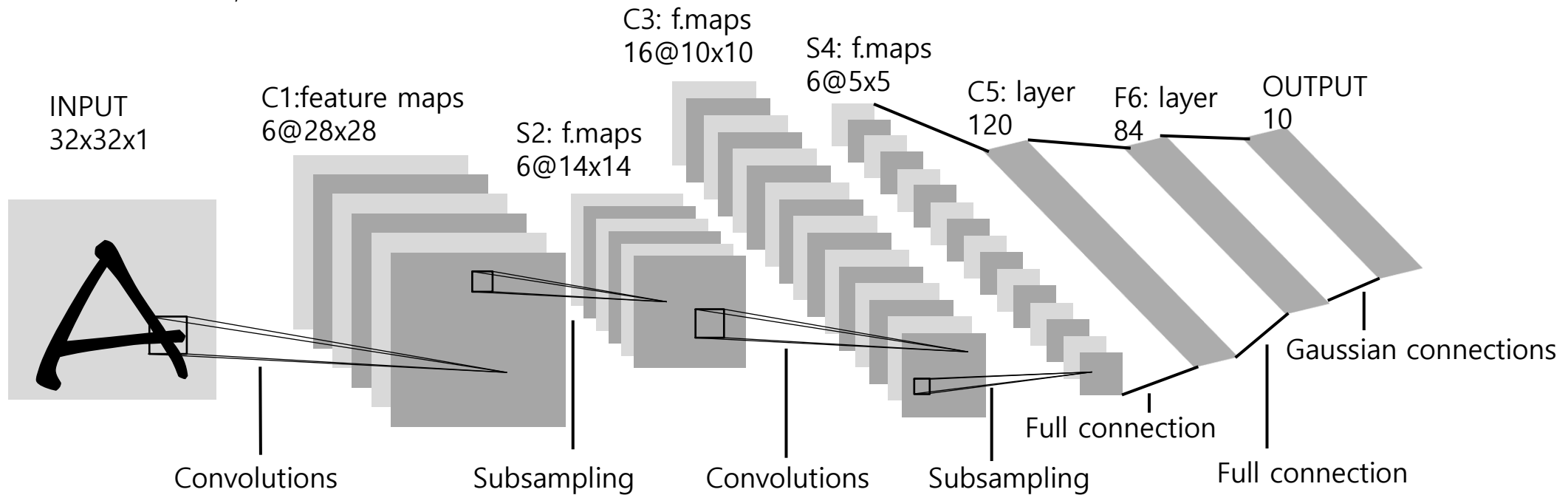
3.5	5.5
11.5	13.5

tf.keras.layers.MAXPool2D

- **pool_size:** integer or tuple of 2 integers, window size over which to take the maximum. (2, 2) will take the max value over a 2x2 pooling window. If only one integer is specified, the same window length will be used for both dimensions.
- **strides:** Integer, tuple of 2 integers, or None. Strides values. Specifies how far the pooling window moves for each pooling step. If None, it will default to pool_size.
- **padding:** One of "valid" or "same" (case-insensitive). "valid" adds no zero padding. "same" adds padding such that if the stride is 1, the output shape is the same as input shape.
- **data_format:** A string, one of channels_last (default) or channels_first. The ordering of the dimensions in the inputs. channels_last corresponds to inputs with shape (batch, height, width, channels) while channels_first corresponds to inputs with shape (batch, channels, height, width). It defaults to the image_data_format value found in your Keras config file at ~/.keras/keras.json. If you never set it, then it will be "channels_last".

LeNet-5

LeCun et al. , 1998



Conv filters 5x5, stride 1
Subsampling 2x2, stride 2

활성화 함수

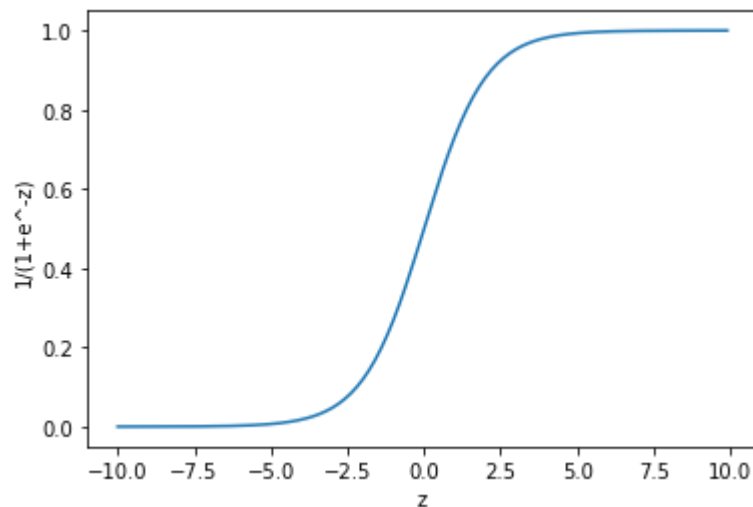
활성함수는 네트워크에 비선형성(nonlinearity)을 추가하기 위해 사용됨

- 활성화 함수 없이 layer를 쌓은 네트워크는 1-layer 네트워크와 동일하기 때문에 활성화 함수는 비선형 함수로 불리기도 한다.
- 멀티레이어 퍼셉트론을 만들 때 활성화 함수를 사용하지 않으면 쓰나마나이다.

1. 시그모이드 함수 (Sigmoid Function)

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

- 결과값이 [0,1] 사이로 제한됨
- 뇌의 뉴런과 유사하여 많이 쓰였음



- 문제점

1) 그레디언트가 죽는 현상이 발생한다 (Gradient vanishing 문제)

gradient 0이 곱해 지니까 그 다음 layer로 전파되지 않는다. 즉, 학습이 되지 않는다.

2) 활성화함수의 결과 값의 중심이 0이 아닌 0.5이다.

3) 계산이 복잡하다 (지수함수 계산)

!! Gradient Vanishing

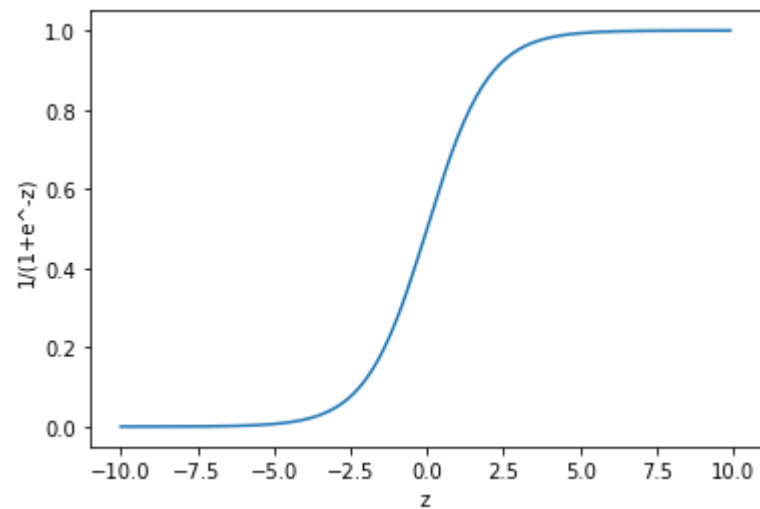
- 시그모이드와 같이 결과값이 포화(saturated)되는 함수는 gradient vanishing 현상을 야기

- 이전 레이어로 전파되는 그라디언트가 0에 가까워 지는 현상

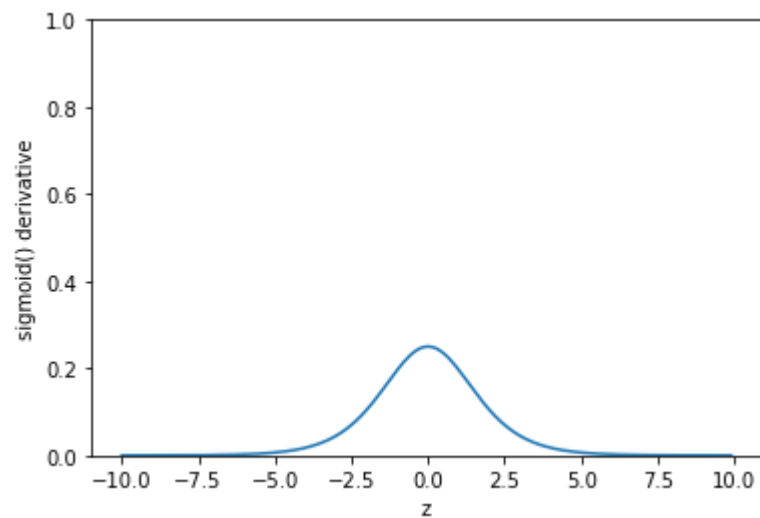
- 레이어를 깊게 쌓으면 파라미터의 업데이트가 제대로 이루어지지 않음

- 양 극단의 미분값이 0에 가깝기 때문에 발생하는 문제

시그모이드 함수



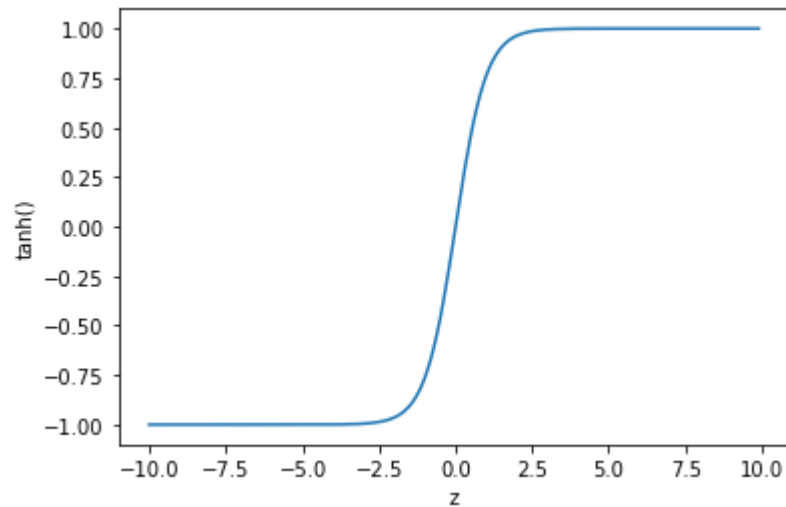
시그모이드 미분값



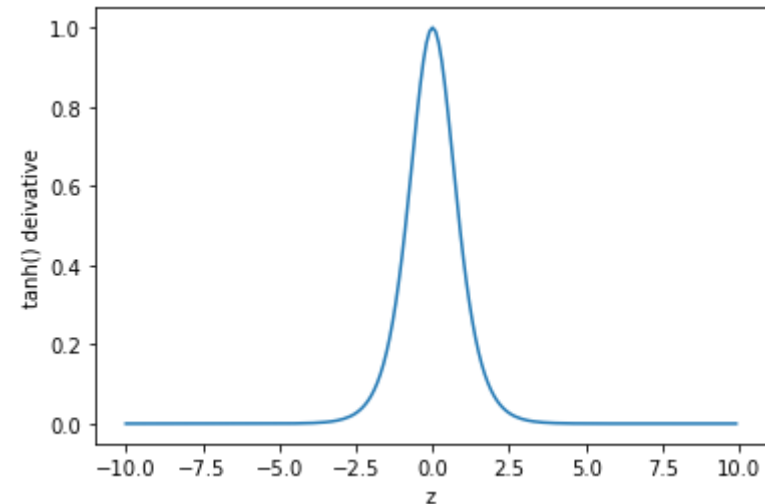
양쪽 꼬리가 0에 수렴하며 최대값이 0.25를 넘지 않는다.

2. 하이퍼 볼릭 탄젠트(tanh)

$$\tanh(x) = 2 * \text{sigmoid}(2 * x) - 1$$



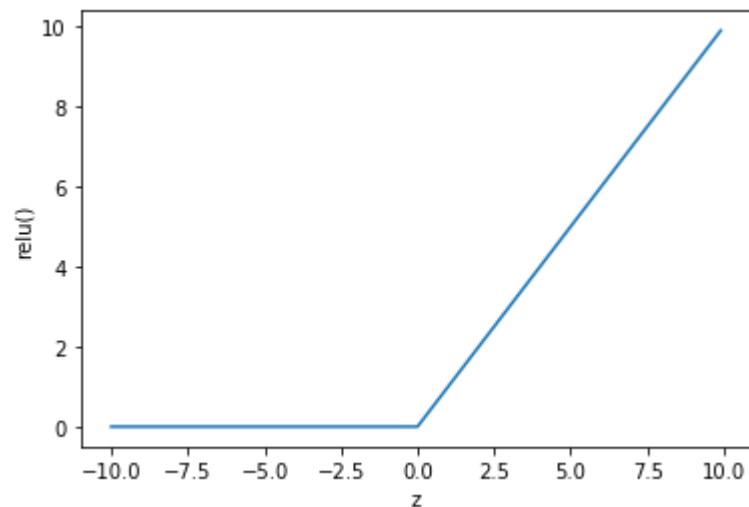
$$\tanh(x)' = (1 - \tanh(x))(1 + \tanh(x))$$



- 결과값이 [-1, 1] 사이로 제한됨. 결과값 중심이 0이다.
- 나머지 특성은 시그모이드와 비슷함. 시그모이드 함수를 이용하여 유도 가능
- 그러나, 여전히 gradient vanishing 문제가 발생

3. 렐루(ReLU, Rectified Linear Unit)

$$f(x) = \max(0, x)$$



최근 뉴럴 네트워크에서 가장 많이 쓰이는 활성화 함수
선형아니야? NO! 0에서 확 꺾이기 때문에 비선형이라고 본다.

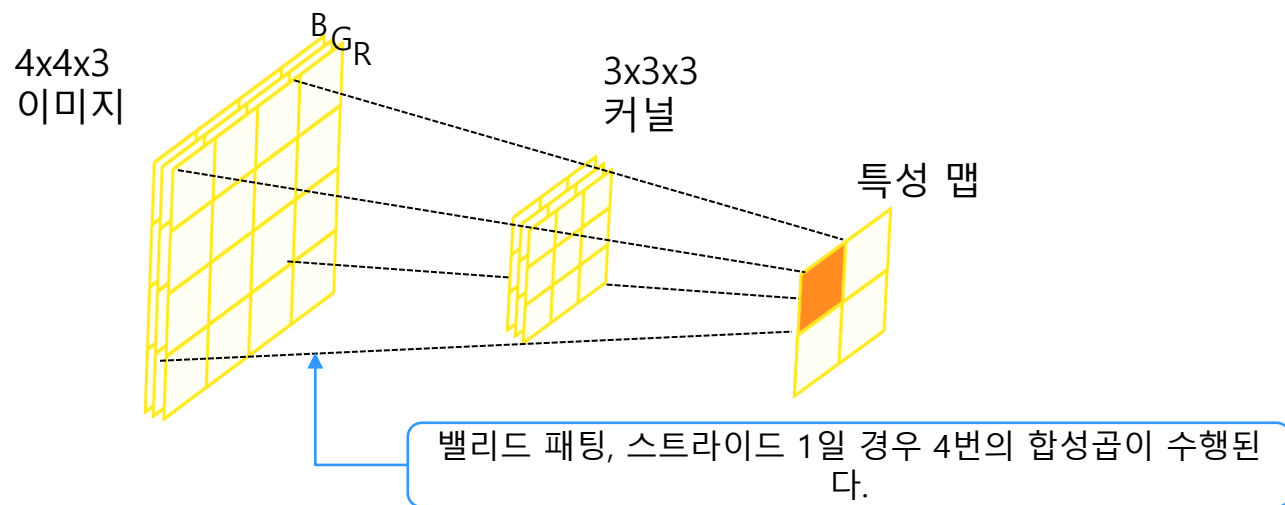
- 장점

- (1) 양 극단값이 포화되지 않는다. (양수 지역은 선형적)
- (2) 계산이 매우 효율적이다 (최대값 연산 1개)
- (3) 수렴속도가 시그모이드류 함수대비 6배 정도 빠르다.

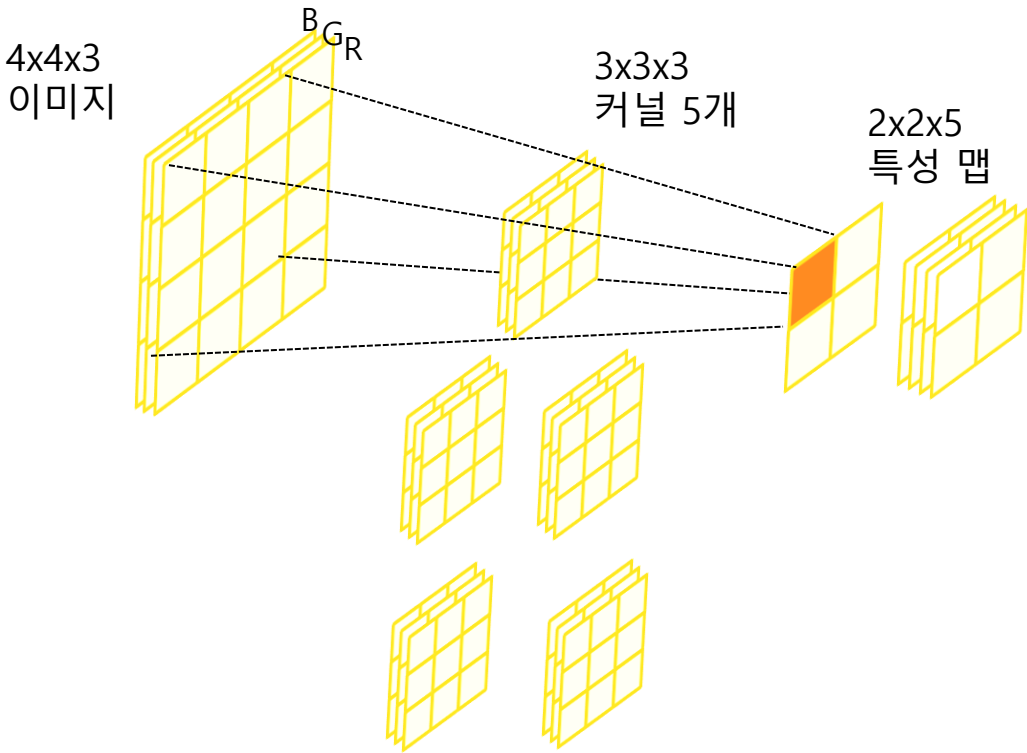
- 단점

- (1) 중심값이 0이 아님 (마이너한 문제)
- (2) 입력값이 음수인 경우 항상 0을 출력함 (마찬가지로 파라미터 업데이트가 안됨)

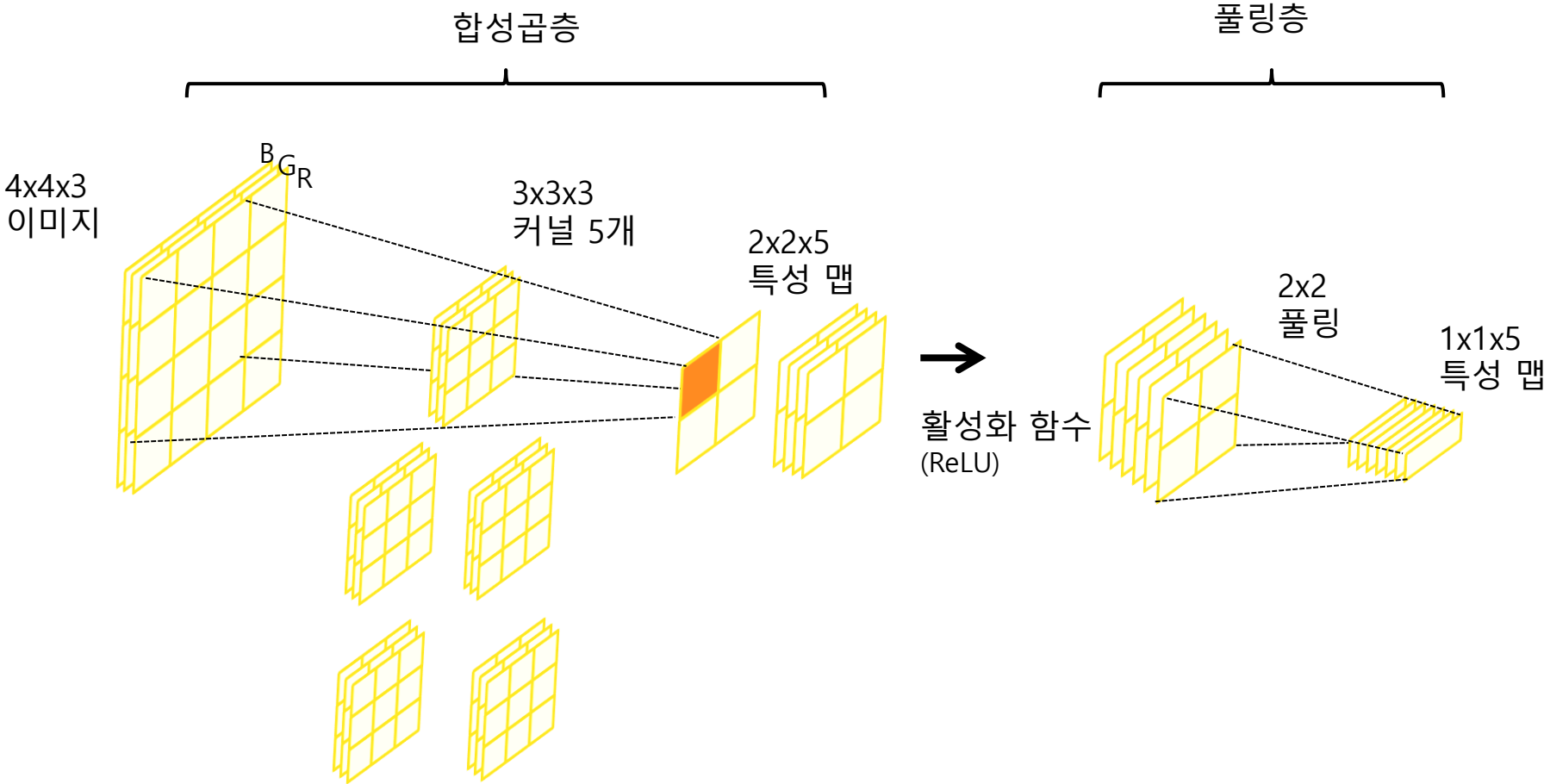
합성곱 층에서 일어나는 일



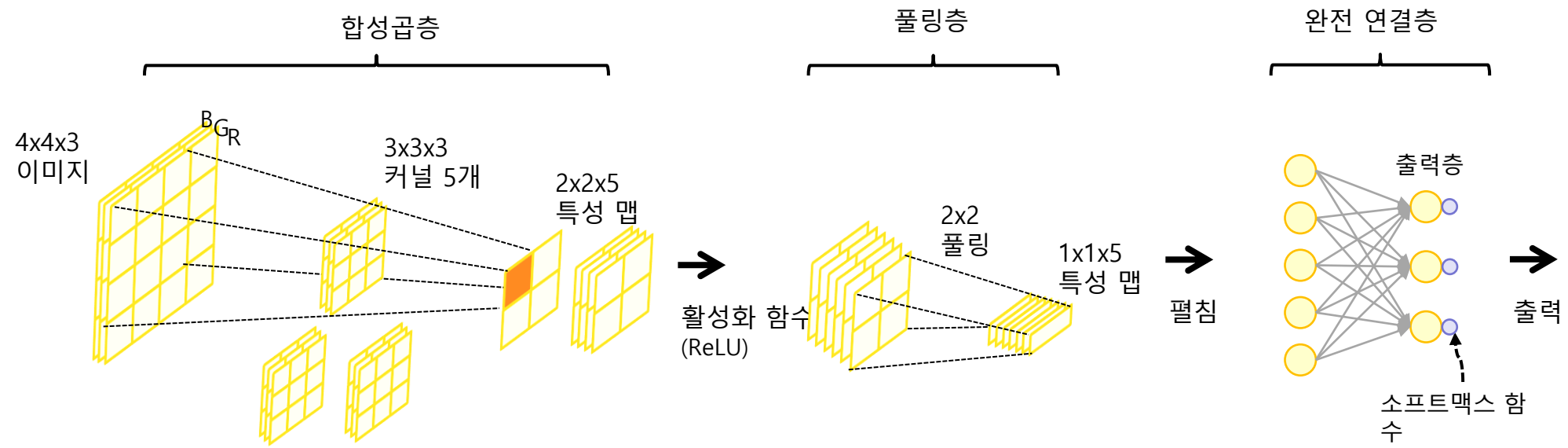
합성곱 층에서 일어나는 일



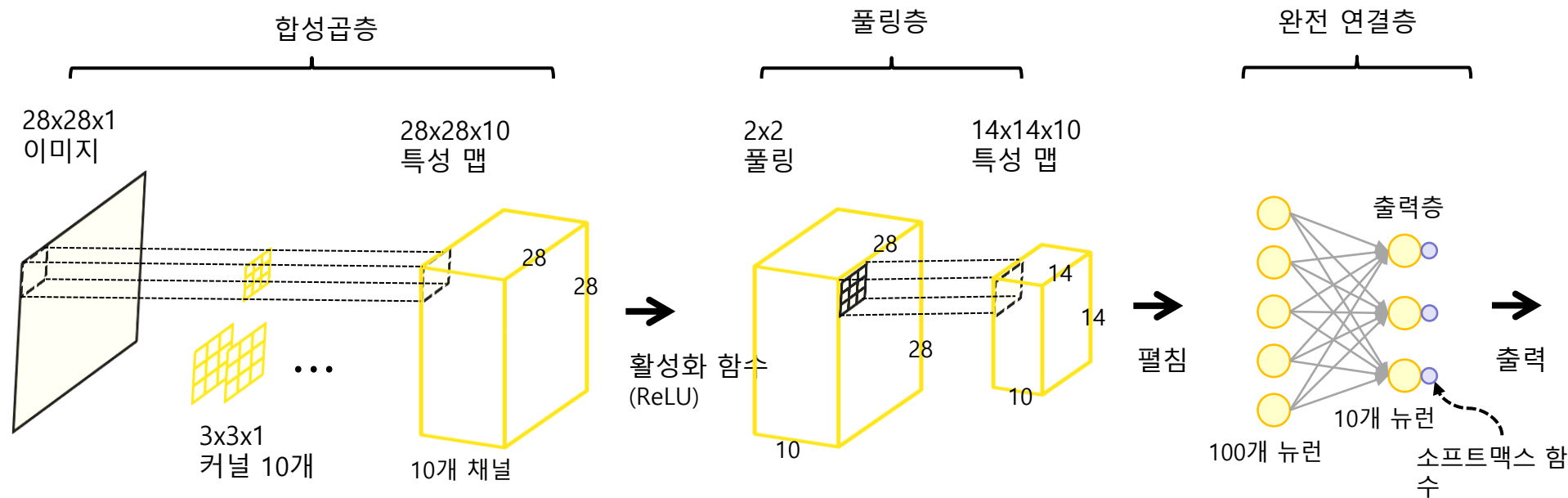
풀링 층에서 일어나는 일



특성 맵을 펼쳐 완전 연결 신경망에 주입한다.



합성곱 신경망의 전체 구조



- 28×28 크기의 흑백 이미지와 3×3 크기의 커널 10개로 합성곱 수행
- 2×2 크기의 최대 풀링을 수행하여 $14 \times 14 \times 10$ 로 특성 맵의 크기를 줄인다.
- 특성 맵을 일렬로 펼쳐서 100개의 뉴런을 가진 완전 연결층과 연결 시킨다.
- 10개의 클래스를 구분하기 위한 소프트맥스 함수에 연결한다.