

Introduction

1) What is Deep Learning?

Artificial Intelligence is mimicing human intelligence

Machine Learning is data-driven approach.

When we want to solve a Problem, We use dataset to find a solution

Deep learning is a part of Machine Learning.

Deep learning uses Neural Networks to solve a problem

Linear regression

Model

- Given a dataset, our goal is to choose the weights w and the bias b that make our model's predictions fit the true price observed in that data as closely as possible.

Input $\rightarrow d$ features. (area와 age로 집가격을 예측한다면
 $\text{Features} = [\text{area}, \text{age}]$)

$$\hat{y} = w_1x_1 + \dots + w_dx_d + b$$

↓
weight

\Rightarrow Affine transformation of input features.

$$\hat{y} = w^T x + b$$

Weight vector ↓ Feature vector ↓

o Weight : a value that adjust the influence of each input feature on the output

o bias : a constant value that adjusts the output of the model, allowing it to shift away from the origin.

Single instance \rightarrow n examples (entire data set)

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

Point (observations)
↓
 $n \times d$ (row)
features (column)

\hookrightarrow Design matrix $X \in \mathbb{R}$

$$\therefore \hat{y} = Xw + b$$

\Rightarrow Broadcasting \rightarrow \hat{y} $\in \mathbb{R}^n$

o Given the features of a training dataset X and corresponding labels y , the goal of linear regression is to find the weight vector w and the bias b with the lowest prediction error.

Model (정리) \rightarrow Loss Function (실제값 y 와 예측값 \hat{y} 의 차의 총합)

; 한개 example이면

$$\text{Squared error } l^{(i)}(w, b) = \frac{1}{2} (\hat{y}^{(i)} - y^{(i)})^2 - \text{미분할 때 계산이 편한데 } \frac{1}{2} \text{ 사용}$$

n example (entire dataset) \rightarrow average (or sum) the losses

$$L(w, b) = \frac{1}{n} \sum_{i=1}^n l^{(i)}(w, b) = \frac{1}{n} \sum_{i=1}^n \frac{1}{2} (\underbrace{\hat{y}^{(i)}}_{\hat{y}} - \underbrace{w^T x^{(i)} + b}_{y^{(i)}})^2$$

loss를 최소화하는 (w^*, b^*) 찾기

$$w^*, b^* = \underset{w, b}{\operatorname{arg\,min}} L(w, b)$$

Analytic Solution

Linear regression에 we can find optimal parameters analytically

$$\text{minimize } \|y - Xw\|^2 \quad \leftarrow \text{square error 가 최소}$$

\leftarrow 전체를 대표하는
(각 예제들의 합)

\downarrow

$$\text{will 대해 연습} \quad \partial_w \|y - Xw\|^2 = 2X^T(Xw - y) = 0$$

\downarrow

$$w^* = (X^T X)^{-1} X^T y$$

w^* 찾기

But, 복잡하기 Analytical 미려을 때!

Stochastic Gradient Descent

Gradient Descent

loss function을 줄이는 방향으로 parameter를 update 하며 error를 줄이는 방향

\rightarrow 같은 방향으로 하면 느림.

\Rightarrow 극단적으로 같은 데이터 헨트만 사용하는게 SGD

\rightarrow 중간 방법으로 minibatch SGD가 있음.

전체 dataset에서 몇개만 선택해 minibatch를 만든다.

$$(w, b) \leftarrow (w, b) - \frac{\eta}{|B|} \sum_{i \in B} \hat{d}_{(w, b)} l^{(i)}(w, b)$$

B : minibatch (샘플의 집합)

$|B|$: minibatch의 크기 (샘플 개수)

η : learning rate (일정한 parameter를 통해 바꿀지)

\therefore SGD iteratively reduces the error by updating parameters using
Only a single example at a time in the direction that lowers the loss function.

Maximum Likelihood Learning

$x \sim N(\mu, \sigma^2)$ mean: μ , variance σ^2

Error가 Normal dist를 따른다 가정.

$$p(y|x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(y-\mu)^2\right)$$

$$y = w^T x + b + \epsilon, \epsilon \sim N(0, \sigma^2)$$

Likelihood function (특정 x 에 대해 y 가 나올 확률)

$$P(y|x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2} \underbrace{(y - (w^T x + b))^2}_{\hat{y} - (w^T x + b)}\right)$$

$\hat{y} \rightarrow E[y|x]$ (MCN 들어감)

The best values of w, b are those that maximize the likelihood of the entire dataset.

For entire dataset,

$$P(y|x) = \prod_{i=1}^n P(y^{(i)}|x^{(i)})$$

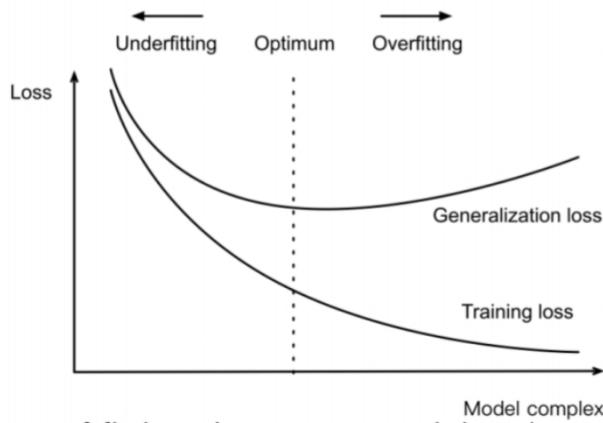
likelihood를 maximize하는 쪽, Negative log-likelihood를 minimize할 수 있다.
(계산 편의)

$$-\log P(y|x) = \sum_{i=1}^n \cancel{\left(\frac{1}{2} \log(2\pi\sigma^2) + \frac{1}{2\sigma^2} (y^{(i)} - w^T x^{(i)} - b)^2 \right)}$$

If σ is fixed, first term 무시 가능

- In fact, as the solution does not depend on σ , minimizing the mean squared error is equivalent to the maximum likelihood estimation of a linear model under the assumption of additive Gaussian noise.

Underfitting or Overfitting



- The phenomenon of fitting closer to our training data than to the underlying distribution of called **overfitting**, and techniques for combatting overfitting are often called **regularization methods**.

o Underfitting

: Underfitting occurs when a model is too simple to capture underlying patterns in the training data. In this case, both the training loss and generalization loss are high, resulting in poor model performance.

o Overfitting

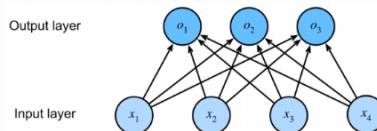
: Overfitting occurs when a model is too complex and fits the training data closely, resulting in poor performance on new data. In this case, the training loss is low, but the generalization loss is high, leading to decreased predictive power.

Classification

Model

ex) $O_{1,2,3} = \text{개}, \text{고양이}, \text{닭}$

$x_{1,2,3,4} = \text{몸무게}, \text{털길이}, \text{날개길이}, \text{귀형태}$



374의 2412213 번호

$$o_1 = x_1 w_{11} + x_2 w_{12} + x_3 w_{13} + x_4 w_{14} + b_1$$

$$o_2 = x_1 w_{21} + x_2 w_{22} + x_3 w_{23} + x_4 w_{24} + b_2$$

$$o_3 = x_1 w_{31} + x_2 w_{32} + x_3 w_{33} + x_4 w_{34} + b_3$$

class \Rightarrow 7 features.

$$\Rightarrow O = Wx + lb, \text{ where } W \in \mathbb{R}^{3 \times 4}, lb \in \mathbb{R}^3$$

Softmax

: In the Classification Setting, Output follows a categorical distribution (i.e., nonnegative and sum up to one.)

=> Softmax utilizes exponential function to model Categorical dist,

$$P(y=i) \propto \exp(o_i)$$

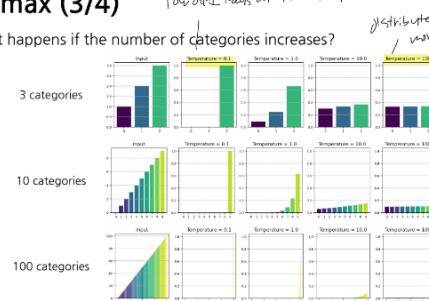
$$\hat{y} = \text{Softmax}(O), \text{ where } \hat{y}_i = \frac{\exp(o_i)}{\sum_j \exp(o_j)}$$

\Rightarrow 여기서 temperature T 를 더해서 Softmax의 output 조절

$$\hat{y} = \text{Softmax}(O), \text{ where } \hat{y}_i = \frac{\exp(o_i/\tau)}{\sum_j \exp(o_j/\tau)}$$

Softmax (3/4)

- What happens if the number of categories increases?

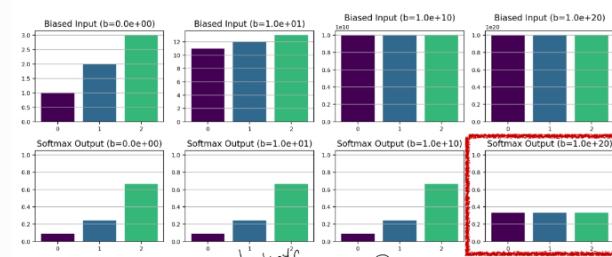


lower T : sharper distribution
(Kelly 243)

higher T : more even distribution

$$\text{biases} \Rightarrow \hat{y} = \text{softmax}(O+b)$$

- What about biases (i.e., $\hat{y} = \text{softmax}(O+b)$ where b is a bias)?



bio는 input이 아니라 출력인거지.

Class \Rightarrow 확률로 바꿔야함.

Vectorization (in minibatches) \Rightarrow 여러 이미지를 한번에
 —> 디비전과 더하기 연산이 없어짐

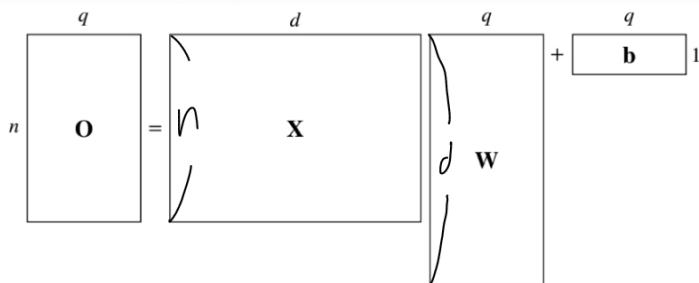
$$\mathcal{O} = XW + b$$

$$\hat{Y} = \text{softmax}(\mathcal{O}) - \text{최종 각 클래스의 확률.}$$

where $\mathcal{O} \in \mathbb{R}^{n \times q}$, $X \in \mathbb{R}^{n \times d}$, $W \in \mathbb{R}^{d \times q}$, $b \in \mathbb{R}^{1 \times q}$, $\hat{Y} \in \mathbb{R}^{n \times q}$

n : 이미지 샘플 개수

q : 클래스 수.



$$\hat{Y} = \begin{pmatrix} \quad & \quad \\ \quad & \quad \\ \quad & \quad \end{pmatrix} \rightarrow \text{이미지 } i \text{의 고양이일 확률, 개일 확률, 새일 확률}$$

Loss Function

Softmax 결과 \hat{Y} : condition probabilities of each class, given input X
 (e.g. $\hat{Y}_i = P(Y=\text{cat} | X)$)

Likelihood

$$P(Y|X) = \prod_{i=1}^n P(Y^{(i)}|X^{(i)})$$

↓ negative log-likelihood

$$-\log P(Y|X) = \sum_{i=1}^n -\log P(Y^{(i)}|X^{(i)}) = \sum_{i=1}^n l(Y^{(i)}|\hat{Y}^{(i)})$$

Where $l(Y, \hat{Y}) = -\sum_{j=1}^q Y_j \log \hat{Y}_j \Rightarrow \text{Cross-entropy loss}$

Cross-Entropy Loss \leftarrow Essential of classification prob

$$l(y, \hat{y}) = -\sum_{j=1}^q y_j \log \hat{y}_j = -\sum_{j=1}^q y_j \log \frac{\exp(o_j)}{\sum_{k=1}^q \exp(o_k)} = \log \frac{1}{\sum_{k=1}^q \exp(o_k)} - \sum_{j=1}^q y_j o_j$$

$$= -\sum_{j=1}^q y_j \log \frac{\exp(o_j)}{\sum_{k=1}^q \exp(o_k)} = \log \frac{\sum_{k=1}^q \exp(o_k)}{\sum_{j=1}^q y_j \exp(o_j)}$$

↓

$$\hat{y} \left(; \text{클래스의 확률} \right) = \frac{\text{클래스의 } \exp \text{ 값}}{\text{모든 클래스의 } \exp \text{ 값}}$$

$$(Se(o) = \log \sum_{k=1}^q \exp(o_k))$$

↓
log-softmax

$$\max(o) \leq Se(o) \leq \max(o) + \log q$$

\Rightarrow give upper bound (적당한 상한 설정)

The derivative of $l(y, \hat{y})$ w.r.t o_j :

$$\partial_{o_j} l(y, \hat{y}) = \frac{\exp(o_j)}{\sum_{k=1}^q \exp(o_k)} - y_j = \text{Softmax}(o_j) - y_j$$

\hat{y}_j

\Rightarrow difference between \hat{y} and y
 ↓
 Probability assigned by model
 label

Information Theory

: to quantify the amount of information contained in data.

For discrete dist P, entropy is

$$H[P] = \sum_j -P(j) \log P(j) = \mathbb{E}[-\log P]$$

확률 분포의 불확실성. -

\Rightarrow expected surprisal (*i.e.*, lower probability, greater surprisal)

negative log Prob (불확실성의 반대)

Cross-entropy from P to Q

$$H(P, Q) = \sum_j -P(j) \log Q(j)$$

Linear model (affine transformation)에서 input과 output의 결과는 linearity라는 특성을 갖.

⇒ But, 실제 문제는 대부분 Non-linear 특성.

⇒ 이를 극복하기 위해, find a suitable representation of inputs,

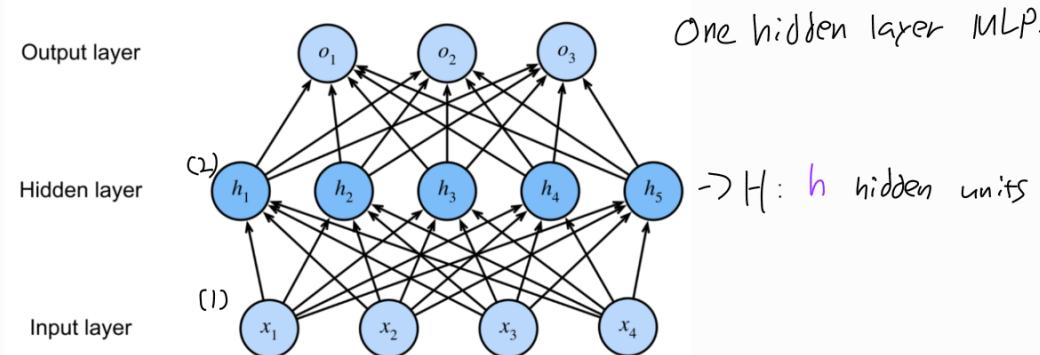
↓

input을 더 적절한 형태로 변형. (원본에 제곱이나 로그 등을 씌어서 새로운 input 만들기)

On top of which a single linear model would suffice

↓

변형된 input의 linear model



$X \in \mathbb{R}^{n \times d}$: minibatch of n examples with d inputs (features)

$H \in \mathbb{R}^{n \times h}$: outputs of the hidden layer. h hidden units

$$\Rightarrow H = XW^{(1)} + b^{(1)}$$

$$O = HW^{(2)} + b^{(2)}$$

However,

$$O = (XW^{(1)} + b^{(1)})W^{(2)} + b^{(2)} = XW + b$$

→ hidden layer를 추가했지만 최종 결과에 변화X.

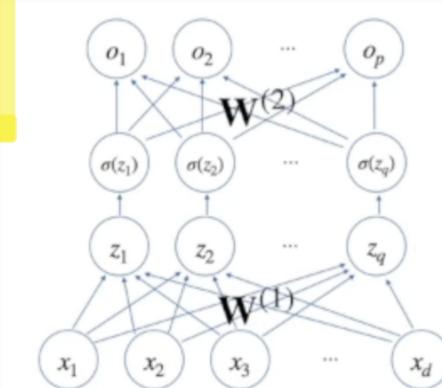
||

Nonlinear activation function σ 필요.

- Popular choice is ReLU (Rectified Linear Unit), $\sigma(x) = \max(0, x)$

$$H = \sigma(XW^{(1)} + b^{(1)})$$

$$O = HW^{(2)} + b^{(2)}$$

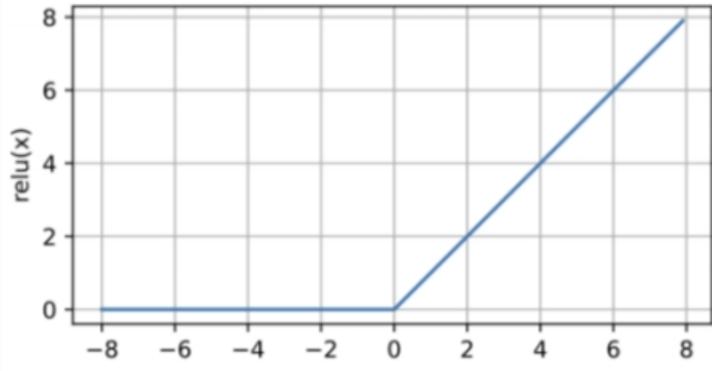


Activation functions.

: differentiable operators, adding non-linearity

여러가지 있는데 각각 장단점이 있음 \Rightarrow Activate function을 고르는 것에 대해 How well network performs on 훈련

ReLU Function: $\text{ReLU}(x) = \max(0, x)$

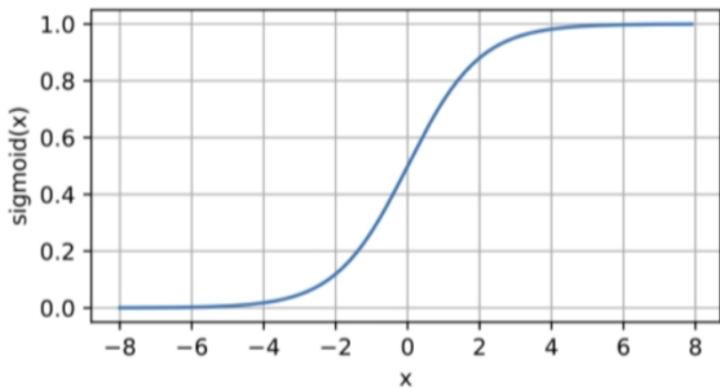


o Vanishing Gradient issue를 다루기 힘들.

* However, "Dead relu" 발생 가능

\Rightarrow Input이 negative면 (bias), large negative result, big output 0.

Sigmoid Function: $\text{Sigmoid}(x) = \frac{1}{1 + \exp(-x)}$



o Output of (0,1) 사이 존재 (확률로 표현 가능)

* - gradient vanishing problem 발생 가능 (input이 너무 작으면 gradient가 매우 작아짐)
(or 너무 크면)

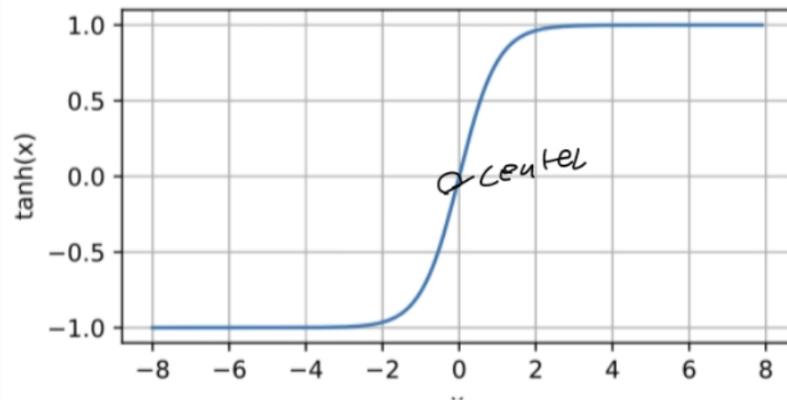
- outputs are not zero-centered.

o useful property

$$\frac{\partial}{\partial x} \text{Sigmoid}(x) = \text{Sigmoid}(x)(1 - \text{Sigmoid}(x))$$

\Rightarrow 계산 편리

Tanh Function: $\tanh(x) = \frac{1 - \exp(-2x)}{1 + \exp(-2x)}$



- Output of $(-1, 1)$ \times 0 \Rightarrow Center 가로 Symmetric

~~Vanish gradient problem~~

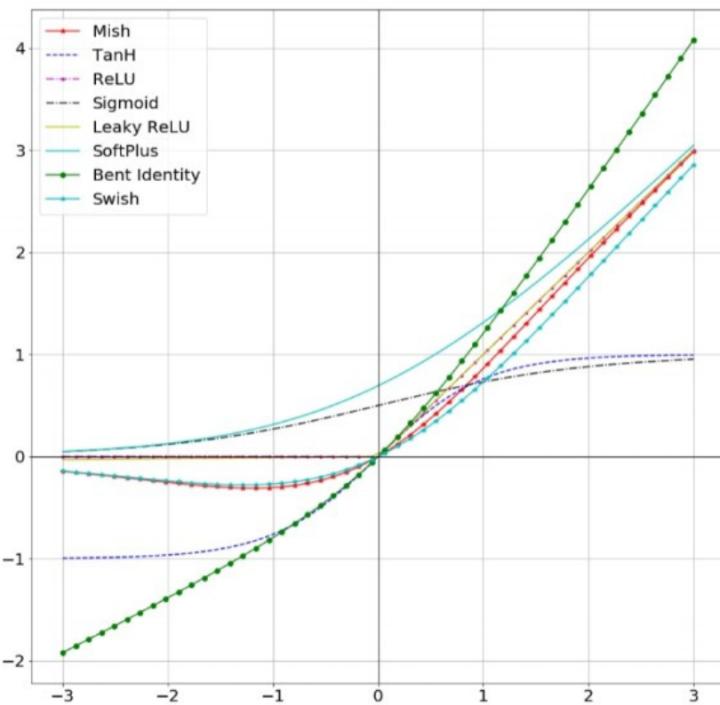
- Useful property

$$\frac{d}{dx} \tanh(x) = 1 - \tanh^2(x)$$

\Rightarrow 계산 쉽다.

2 21

Activation Functions



• **relu** is notorious for the dead relu problem.
• **elu** is proposed to handle this.

Activation value is always zero



- The **relu** is notorious for the dead relu problem.
- To handle this, the **elu** function was proposed. However, it introduces a longer computation time due to the exponential operation included.

$$\text{elu}(x) = \begin{cases} x, & \text{if } x > 0 \\ \alpha(e^x - 1), & \text{if } x < 0 \end{cases}$$

- The **leaky relu** function also avoids the dead relu problem and is fast. However, we have to tune the parameter α .

$$\text{lrelu}(x) = \begin{cases} x, & \text{if } x > 0 \\ \alpha x, & \text{if } x < 0 \end{cases}$$

- The **gelu** function works well in NLP, specifically **Transformer** models, as it is fast.

$$\text{gelu}(x) = 0.5x(1 + \tanh(\sqrt{2/\pi}(x + 0.044715x^3)))$$

- The **swish** function is continuous and differentiable at all points. And it works well on standard image datasets (CIFAR or ImageNet) compared to others (relu, lrelu, elu, gelu).

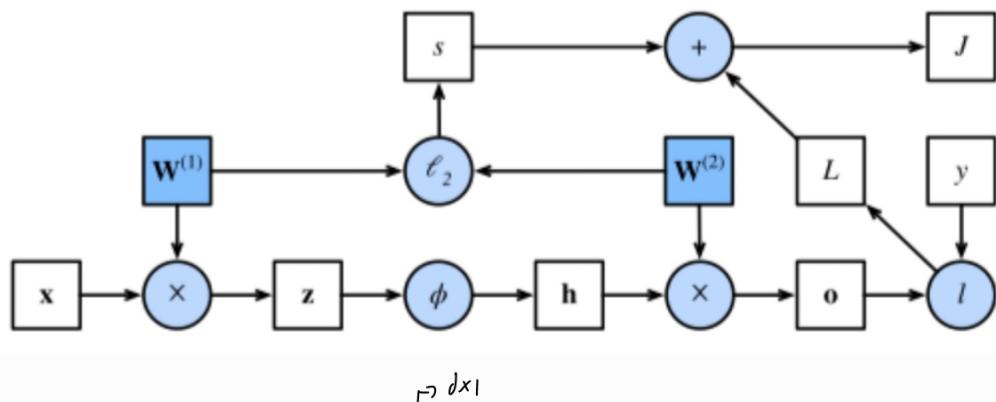
$$\text{swish}(x) = x(1 + e^{-x})^{-1}$$

- The **mish** function is C^∞ -continuous and approximates identity near the origin. In some experiments, the **mish** works better than **swish** activations.

$$\text{mish}(x) = x \tanh(\text{softplus}(x))$$

Back Propagation

Forward Propagation



1) Input: $\mathbf{x} \in \mathbb{R}^d$ (no bias for simplicity)

2) Intermediate Variable: $\mathbf{z} = \mathbf{W}^{(1)}\mathbf{x} \in \mathbb{R}^h$, where $\mathbf{W}^{(1)} \in \mathbb{R}^{h \times d}$ (weight parameter of the hidden layer)

3) Hidden activation Vector: $\mathbf{h} = \phi(\mathbf{z}) \in \mathbb{R}^h \rightarrow$ hidden layer
 ↳ activation function

4) Hidden layer output (MLP 2nd output)

; $\mathbf{o} = \mathbf{W}^{(2)}\mathbf{h} \in \mathbb{R}^q$, where $\mathbf{W}^{(2)} \in \mathbb{R}^{q \times h}$

(weight parameter of the hidden layer)

5) Loss term (for a single data): $L = \int_{-\text{loss function}}^{\text{loss function}} (\mathbf{o}, \mathbf{y})$

, 정규화의 정도

6) Regulation term : $S = \frac{\lambda}{2} (\|\mathbf{W}^{(1)}\|_F^2 + \|\mathbf{W}^{(2)}\|_F^2)$

⇒ Prevent over fitting

$\|\mathbf{W}\|_F \Rightarrow$ Frobenius Norm

$$\sqrt{w_1^2 + w_2^2 + w_3^2 + \dots}$$

7) Model's regularized loss: $J = L + S$

⇒ $J \equiv$ Minimize $\partial L + \lambda \|\mathbf{W}\|_F^2$

Back propagation

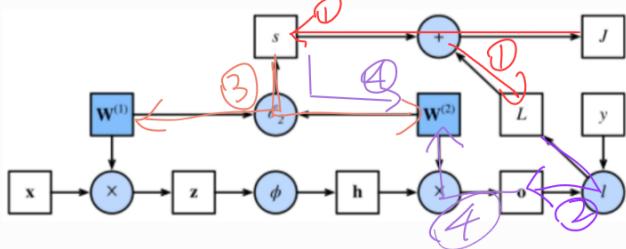
: Calculating the gradient of neural network parameters
(to update weights)

-> Chain rule을 사용해 출력에서 input으로 향해

Chain rule

$$Y = f(x), Z = g(Y)$$

$$\frac{\partial Z}{\partial x} = \prod \left(\frac{\partial Z}{\partial Y}, \frac{\partial Y}{\partial x} \right)$$



$$\mathbf{o} = \mathbf{W}^{(2)} \mathbf{h} \in \mathbb{R}^q$$

$$s = \frac{\lambda}{2} (\|\mathbf{W}^{(1)}\|_F^2 + \|\mathbf{W}^{(2)}\|_F^2)$$

$\mathbf{W}^{(2)}$

1) Calculate the gradients of the objective $J = L + s$ w.r.t $\frac{\partial L}{\partial z}$ and $\frac{\partial s}{\partial z}$.

Loss term regularization term

$$\frac{\partial J}{\partial L} = 1, \quad \frac{\partial J}{\partial s} = 1$$

$$\Leftrightarrow \frac{\partial (L + s)}{\partial L} = 1$$

2) Compute the gradient of J w.r.t the output layer \mathbf{o} : \mathbf{o} 가 J 에 얼마나 영향을 미치는지?

$$\frac{\partial J}{\partial \mathbf{o}} = \prod \left(\frac{\partial J}{\partial L}, \frac{\partial L}{\partial \mathbf{o}} \right) = \frac{\partial L}{\partial \mathbf{o}} \in \mathbb{R}^q$$

각 클래스별 해석
 $\Rightarrow \mathbb{R}^q$

$$\prod \left(\frac{\partial J}{\partial L}, \frac{\partial L}{\partial \mathbf{o}} \right) = \frac{\partial J}{\partial L} \times \frac{\partial L}{\partial \mathbf{o}}$$

\Rightarrow chain rule

3) Calculate the gradients of s w.r.t Parameters $\mathbf{W}^{(1)}, \mathbf{W}^{(2)}$

$$\frac{\partial s}{\partial \mathbf{w}^{(1)}} = \lambda \mathbf{w}^{(1)}, \quad \frac{\partial s}{\partial \mathbf{w}^{(2)}} = \lambda \mathbf{w}^{(2)}$$

$$\Leftrightarrow s = \frac{\lambda}{2} (\|\mathbf{w}^{(1)}\|_F^2 + \|\mathbf{w}^{(2)}\|_F^2)$$

↑
모든 계층에 대한 정제

$$\frac{\partial s}{\partial \mathbf{w}^{(1)}} = \frac{\lambda}{2} \cdot 2 \mathbf{w}^{(1)} + 0$$

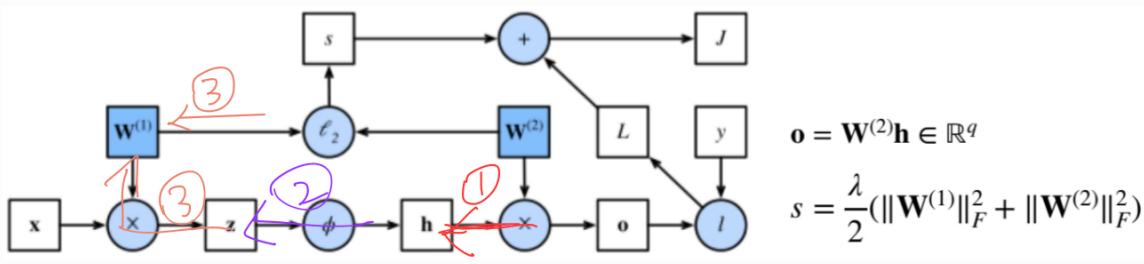
4) Calculate the gradient $\frac{\partial J}{\partial \mathbf{w}^{(2)}} \in \mathbb{R}^{q \times h}$

근원에서 차운 열차위해

$$\frac{\partial J}{\partial \mathbf{w}^{(2)}} = \prod \left(\frac{\partial J}{\partial \mathbf{o}}, \frac{\partial \mathbf{o}}{\partial \mathbf{w}^{(2)}} \right) + \prod \left(\frac{\partial J}{\partial s}, \frac{\partial s}{\partial \mathbf{w}^{(2)}} \right) = \frac{\partial J}{\partial \mathbf{o}} \mathbf{h}^T + \lambda \mathbf{w}^{(2)}$$

$\mathbf{o} = \mathbf{W}^{(2)} \mathbf{h}$

$$\frac{\partial \mathbf{o}}{\partial \mathbf{w}^{(2)}} = \mathbf{h}$$



$\mathbf{W}^{(1)}$

1) To obtain the gradient w.r.t $\mathbf{W}^{(1)}$,

We need to Continue back propagation along the output layer to hidden layer

$$\frac{\partial J}{\partial \mathbf{h}} = \prod \left(\frac{\partial J}{\partial \mathbf{o}}, \frac{\partial \mathbf{o}}{\partial \mathbf{h}} \right) = \mathbf{W}^{(2)\top} \cdot \underbrace{\frac{\partial J}{\partial \mathbf{o}}}_{\mathbf{Q} \times \mathbf{1}} \quad \downarrow$$

$$\frac{\partial \mathbf{o}}{\partial \mathbf{h}} = \mathbf{W}^{(1)}$$

2) Since the activation function ϕ applies element wise.

$$\frac{\partial J}{\partial \mathbf{z}} = \prod \left(\frac{\partial J}{\partial \mathbf{h}}, \frac{\partial \mathbf{h}}{\partial \mathbf{z}} \right) = \frac{\partial J}{\partial \mathbf{h}} \odot \phi'(\mathbf{z})$$

\uparrow 逐元素
 $\mathbf{h} = \phi(\mathbf{z})$
 \uparrow 元素级
 $\begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix}, \begin{bmatrix} e \\ f \\ g \\ h \end{bmatrix}$
 $a \odot b = \begin{bmatrix} a \times e \\ b \times f \\ c \times g \\ d \times h \end{bmatrix}$

3) the gradient $\frac{\partial J}{\partial \mathbf{w}^{(1)}} \in \mathbb{R}^{h \times d}$

$$\frac{\partial J}{\partial \mathbf{w}^{(1)}} = \prod \left(\frac{\partial J}{\partial \mathbf{z}}, \frac{\partial \mathbf{z}}{\partial \mathbf{w}^{(1)}} \right) + \prod \left(\frac{\partial J}{\partial \mathbf{s}}, \frac{\partial \mathbf{s}}{\partial \mathbf{w}^{(1)}} \right) = \frac{\partial J}{\partial \mathbf{z}} \star^T + \lambda \mathbf{w}^{(1)}$$

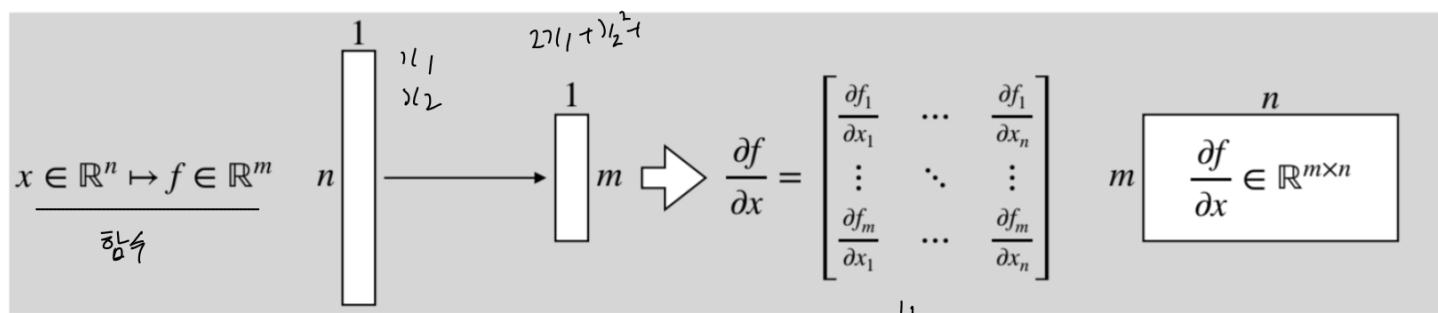
\uparrow 逐元素
 $\mathbf{z} = \mathbf{w}^{(1)} \star$

$$\frac{\partial \mathbf{z}}{\partial \mathbf{w}^{(1)}} = \mathbf{x}$$

$$\begin{aligned}
 \mathbf{w}^{(2)\top} \frac{\partial J}{\partial \mathbf{o}} &= \left(\frac{\partial J}{\partial \mathbf{h}} \odot \phi'(\mathbf{z}) \right) \star^T + \lambda \mathbf{w}^{(2)} \\
 &= \left((\mathbf{w}^{(2)\top} \frac{\partial L}{\partial \mathbf{o}}) \odot \phi'(\mathbf{w}^{(1)} \star) \right) \star^T + \lambda \mathbf{w}^{(2)}
 \end{aligned}$$

Matrix Calculus (행렬 미적분학)

$$\begin{bmatrix} & \\ & \end{bmatrix}$$



Ex)

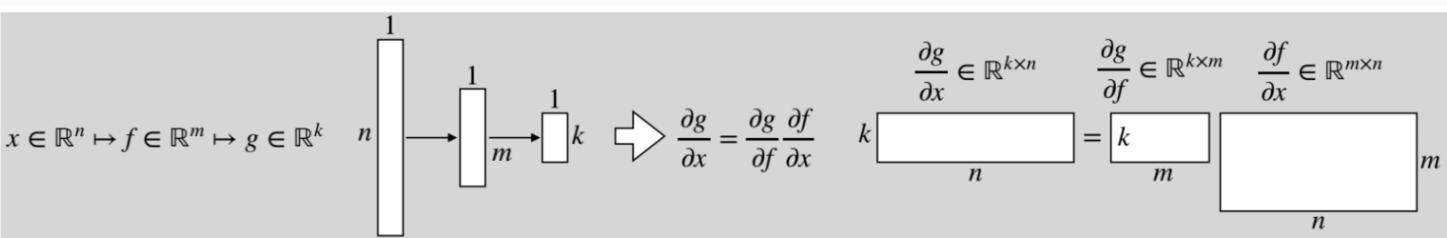
$$x \in \mathbb{R}^2, \quad x = [x_1, x_2]^T$$

Jacobian Matrix

$$f \in \mathbb{R}^3, \quad f(x) = [2x_1 + x_2^2, -x_1^2 + 3x_2, 4x_1 x_2]^T$$

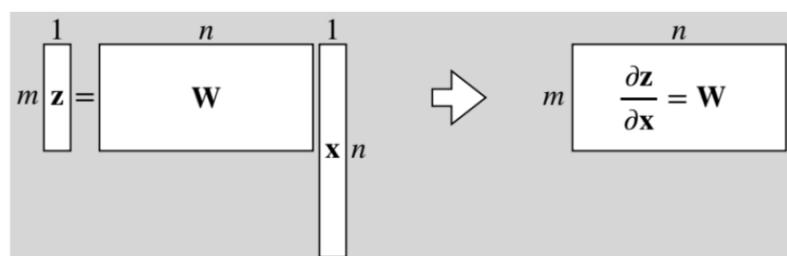
$$\frac{\partial f}{\partial x} = \begin{bmatrix} 2 & 2x_2 \\ -2x_1 & 3 \\ 4x_2 & 4x_1 \end{bmatrix}$$

$$\left(\begin{array}{c|c} x_1 & \\ \hline x_2 & \end{array} \right) \quad \left(\begin{array}{c|c} 2x_1 + x_2^2 & \\ \hline -x_1^2 + 3x_2 & \\ \hline 4x_1 x_2 & \end{array} \right)$$

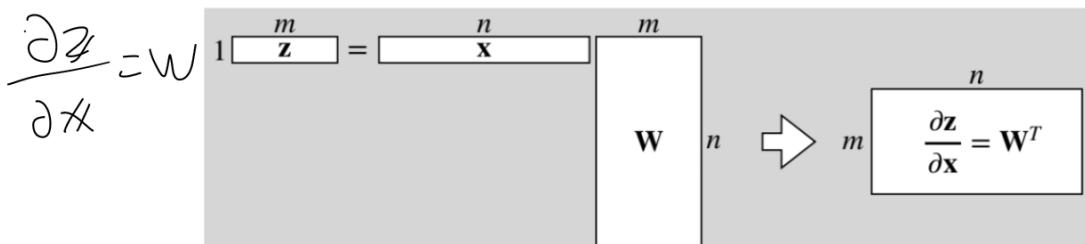


Chain rule

- Matrix times column vector: $\mathbf{z} = \mathbf{W}\mathbf{x}$ where $\mathbf{z} \in \mathbb{R}^{m \times 1}$, $\mathbf{x} \in \mathbb{R}^{n \times 1}$, and $\mathbf{W} \in \mathbb{R}^{m \times n}$. Then, $\frac{\partial \mathbf{z}}{\partial \mathbf{x}} = \mathbf{W} \in \mathbb{R}^{m \times n}$.



- Row vector times matrix: $\mathbf{z} = \mathbf{x}\mathbf{W}$ where $\mathbf{z} \in \mathbb{R}^{1 \times m}$, $\mathbf{x} \in \mathbb{R}^{1 \times n}$, and $\mathbf{W} \in \mathbb{R}^{n \times m}$. Then, $\frac{\partial \mathbf{z}}{\partial \mathbf{x}} = \mathbf{W}^T \in \mathbb{R}^{m \times n}$.



$$1) \quad Z = X. \frac{\partial Z}{\partial X} = I. \begin{bmatrix} 1 & & \\ & \ddots & \\ & & 1 \end{bmatrix}$$

↓
identity matrix

2) An elementwise function applied to a vector $Z = f(X)$, where $Z_i = f(x_i)$

$$\frac{\partial Z}{\partial X} = \text{diag}(f'(X)) \quad \text{or} \quad \odot f'(X)$$

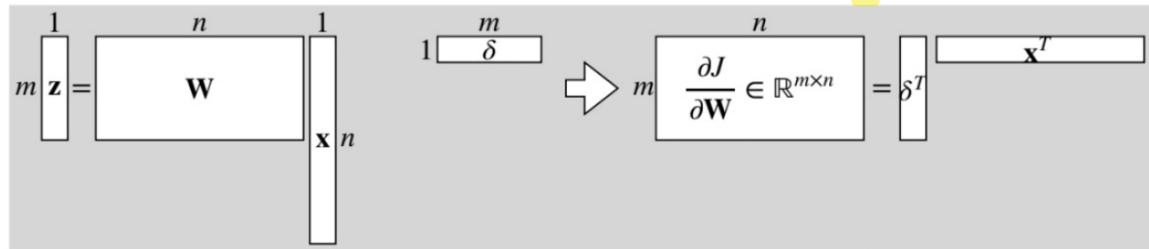
$$Ex) \quad X = [x_1, x_2, x_3]^T$$

$$f(X) = [x_1^2, x_2^2, x_3^2]^T$$

$$\frac{\partial Z}{\partial X} = \begin{bmatrix} 2x_1 & & \\ & 2x_2 & \\ & & 2x_3 \end{bmatrix}$$

gradient of loss function J w.r.t Z

• Matrix times column vector: $z = Wx$ and $\delta = \frac{\partial J}{\partial z} \in \mathbb{R}^{1 \times m}$. Then, $\frac{\partial J}{\partial W} = \frac{\partial J}{\partial z} \frac{\partial z}{\partial W} = \delta \frac{\partial z}{\partial W} = \delta^T x^T \in \mathbb{R}^{m \times n}$



$$X : n \times 1$$

$$\delta^T \cdot X^T = (m \times 1) \cdot (1 \times n) = m \times n$$

• Row vector times Matrix: $x \in \mathbb{R}^{1 \times n}$, $W \in \mathbb{R}^{n \times m}$, $z = xW \in \mathbb{R}^{1 \times m}$, and $\delta = \frac{\partial J}{\partial z} \in \mathbb{R}^{1 \times m}$.

$$\text{Then, } \frac{\partial J}{\partial W} = \frac{\partial J}{\partial z} \frac{\partial z}{\partial W} = \delta \frac{\partial z}{\partial W} = \delta^T \delta \in \mathbb{R}^{n \times m}$$

2:

$$n \times 1 \quad 1 \times m$$

Cross-entropy loss w.r.t. logits

↳ Multi classification problem.

$$J = C \in (\gamma, \hat{\gamma}) . \quad \hat{\gamma} = \underbrace{\text{softmax}(z)}_{\text{logit (raw output)}} \in \mathbb{R}^{l \times k}$$

$$\frac{\partial J}{\partial z} = \hat{\gamma} - \gamma \in \mathbb{R}^{l \times h} . \quad l \text{ is the number of classes}$$

↓

gradient ∂

Vanishing and Exploding Gradients.

o L layers, input x , output θ

o each layer l defined by f_l parametrized by $W^{(l)}$,

whose hidden layer output $h^{(l)}$ (e.g. $h^{(0)} = x$)

$$h^{(l)} = f_l(h^{(l-1)}) , \quad \theta = f_L \circ \dots \circ f_1(x)$$

↓

$$\text{ex) } h^{(1)} = f_1(x)$$

$$h^{(2)} = f_2(h^{(1)}) = f_2(f_1(x))$$

o gradient of θ w.r.t. $W^{(l)}$

$$\partial_{W^{(l)}} \theta = \partial_{h^{(l-1)}} h^{(l)} \dots \partial_{h^1} h^{(l+1)} \partial_{W^{(l)}} h^{(l)}$$

$$\frac{\partial \theta}{\partial W^{(l)}}$$

W 가 변했을 때 θ 가 어떻게 변화하는지

Other Issues

Nonparametrics. ppt 2

MLP \Rightarrow parametric (have fixed # of params)

Nonparametric model

: a level of complexity grows as the amount of available data grows.

\Rightarrow Because neural networks are over-parametrized,
they tend to interpolate the training data (fitting it perfectly) having the same property as nonparametric models.

Regularization in Neural Networks.

(Prevent overfitting) ppt 2

{ Early stopping
weight decay
Dropout .

Invariance

- Despite Waldo's characteristic outfit, what Waldo looks like **does not depend upon where is located.** 위치에 상관X
- We could sweep the image with a Waldo detector that could assign a score to each patch, indicating **the likelihood that the patch contains Waldo.**
- Convolutional neural networks (CNNs) systematize this idea of **spatial invariance**, exploiting it to learn useful representations with **fewer parameters.** → more effective to image prob

No matter where an object appears in image, its feature remain consistent

Convolutions (in single processing)

Convolution between 2 functions, f and g

$$(f * g)(x) = \int f(z)g(x-z)dz$$

$\Rightarrow g$ 함수를 flip, x 만큼 shift.

o Discrete objects in 2-dimensional spaces.

$$(f * g)(i, j) = \sum_a \sum_b f(a, b)g(i-a, j-b) \Rightarrow \text{이미지 처리의 표현.}$$

$f(a, b) \Rightarrow \text{이미지의 특성 필터 (자장)}$

$g(i-a, j-b) \Rightarrow \text{적용할 필터의 값}$

Convolutions (in CNNs)

To exploit the translation invariance, an output of a convolutional layer is computed as:

$$[H]_{i,j} = \underbrace{U}_{\substack{\text{Output} \\ \text{Pixel Position}}} + \sum_{a=-\Delta}^{\Delta} \sum_{b=-\Delta}^{\Delta} [V]_{a,b} [X]_{i+a, j+b}$$

↑
bias
↑
size of filter.
↑
Convolution filter (kernel)
↑
Input

The number of parameters for the layer: $(2\Delta+1)^2$

Convolution

Kernel \leftarrow

Input

3	3	2	1	0
0	0	1	2	3
3	2	1	2	3
2	0	0	2	2
2	0	0	0	1

Kernel

3	3	2	1	0
0	0	1	2	3
3	2	1	2	3
2	0	0	2	2
2	0	0	0	1

Output

3	3	2	1	0
12.0	12.0	17.0		
10.0	17.0	19.0		
9.0	6.0	14.0		

3	3	2	1	0
0	0	1	2	3
3	2	1	2	3
2	0	0	2	2
2	0	0	0	1

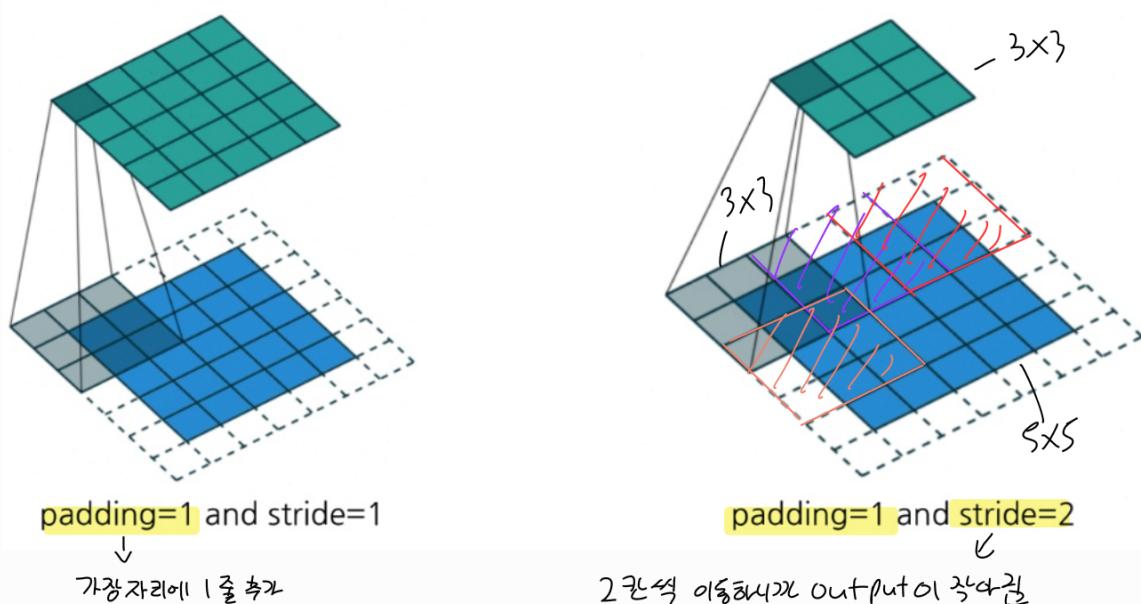
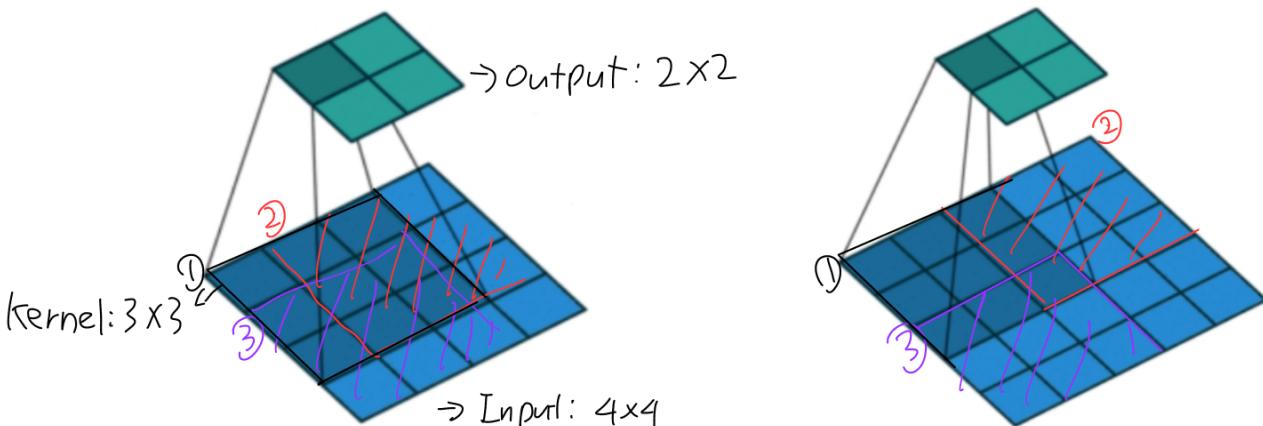
3	3	2	1	0
0	0	1	2	3
3	2	1	2	3
2	0	0	2	2
2	0	0	0	1

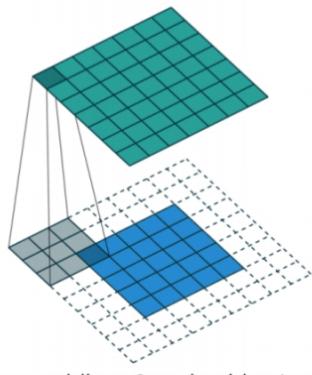
3	3	2	1	0
0	0	1	2	3
3	2	1	2	3
2	0	0	2	2
2	0	0	0	1

3	3	2	1	0
0	0	1	2	3
3	2	1	2	3
2	0	0	2	2
2	0	0	0	1

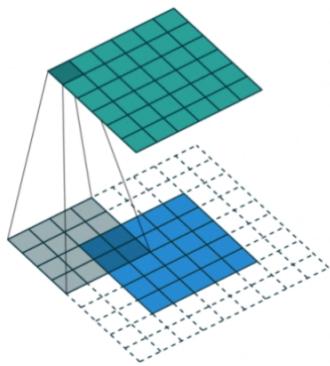
3	3	2	1	0
0	0	1	2	3
3	2	1	2	3
2	0	0	2	2
2	0	0	0	1

Convolution (Padding and Strides)





padding=2 and stride=1



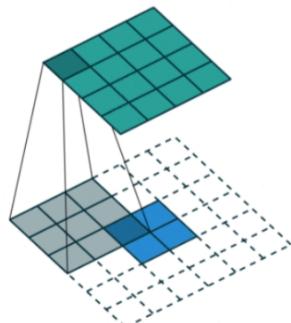
padding=2 and stride=1

Note that the sizes of the output features vary!

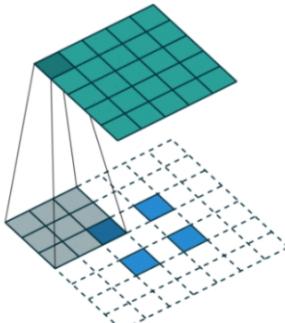


Transposed Convolution (Padding and Strides)

↳ 작은 input으로 큰 output (upsampling)

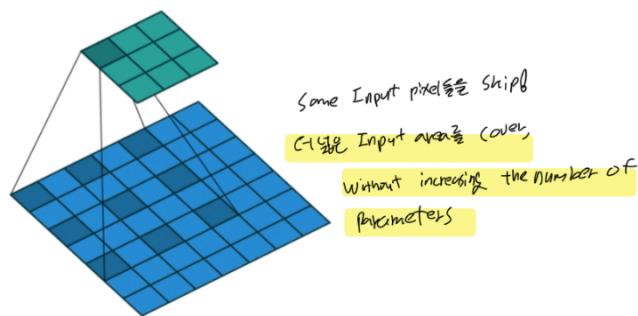


padding=0 and stride=0



padding=0 and stride=1

Dilated Convolution (Padding and Strides)



padding=0 and stride=1

Pooling

way of reducing dimension

Input

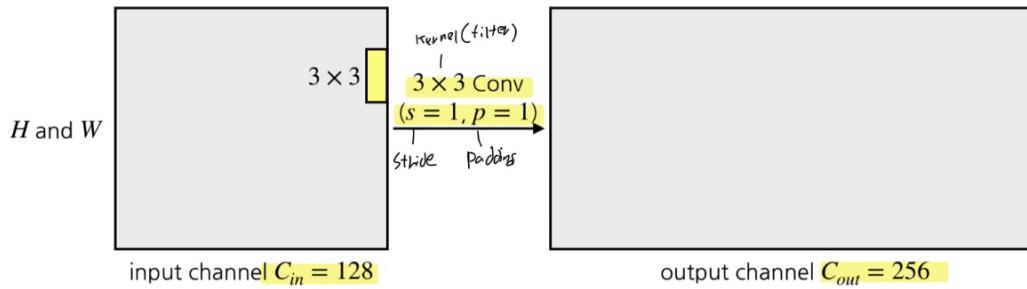
0	1	2
3	(4)	(5)
6	(7)	(8)

Output

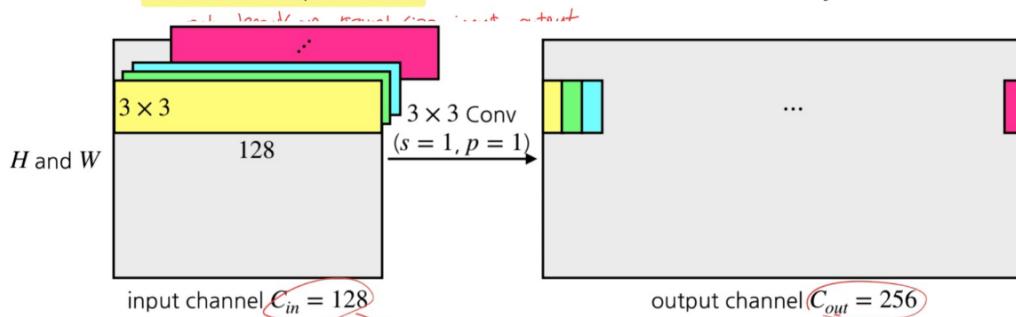
2 x 2 Max Pooling

4	5
7	8

Channel



- What is the number of parameters of this 3×3 convolution layer?



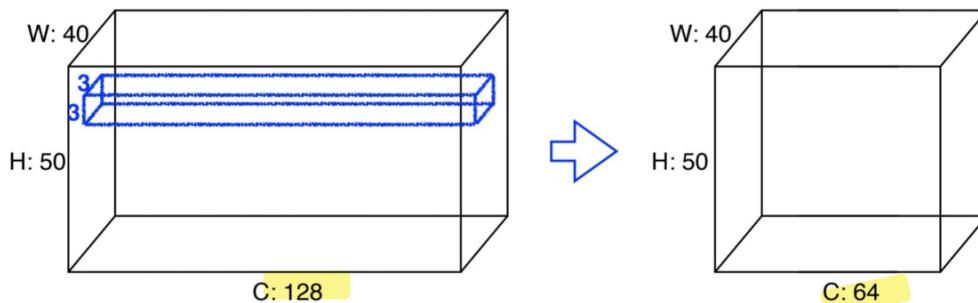
- What is the number of parameters of this 3×3 convolution layer?

- Each kernel has a dimension of $128 \times 3 \times 3$
- And we have 256 kernels.
- The number of parameters is $128 \times 256 \times 3 \times 3 = 294,912$

17

$$3 \times 3 \times 128 \times 256$$

- Padding (1), Stride (1), 3×3 Kernel



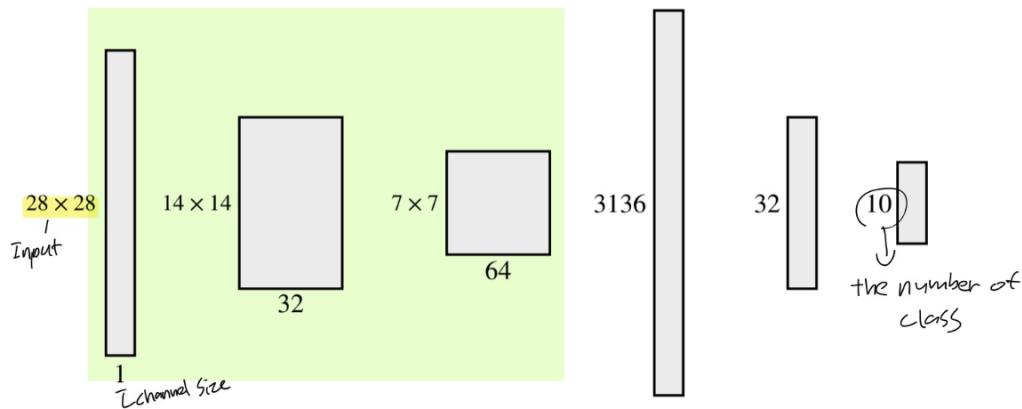
What is the number of parameters of this model? $P, S, W, H \rightarrow$ 상관 x

The answer is $3 \times 3 \times 128 \times 64 = 73,728$

19

$$128 \times 64 \times 3 \times 3$$

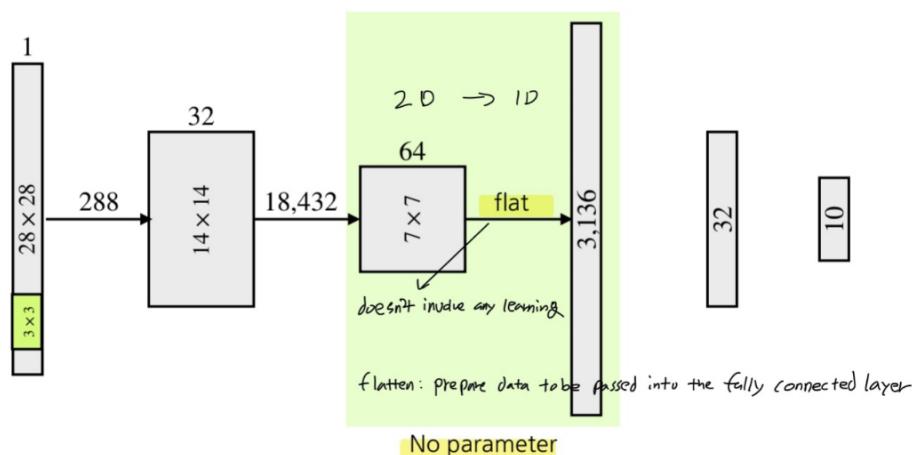
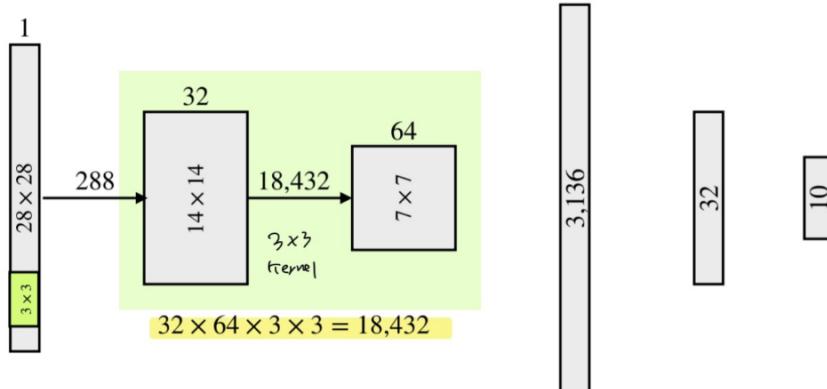
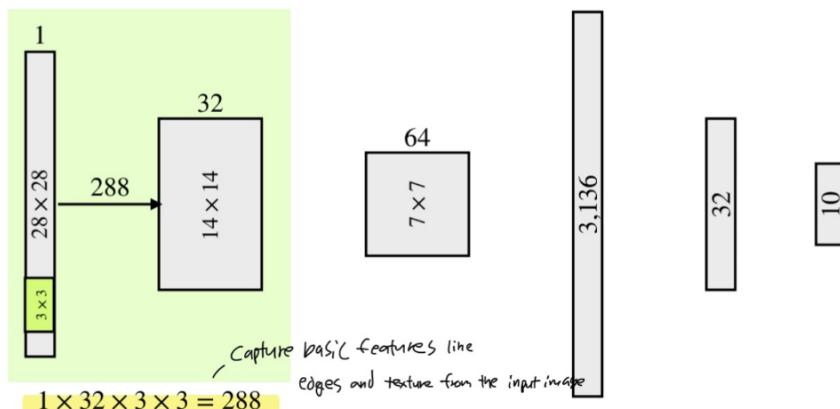
Convolutional Neural Network



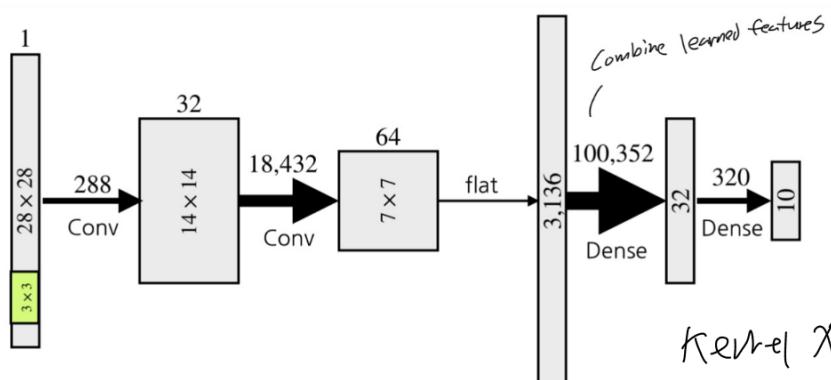
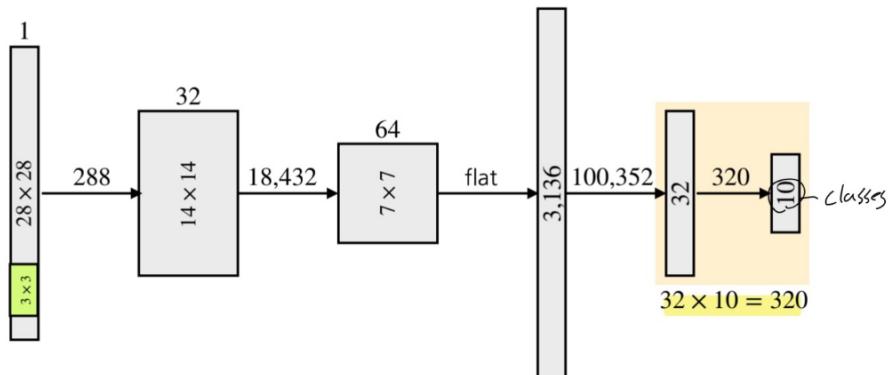
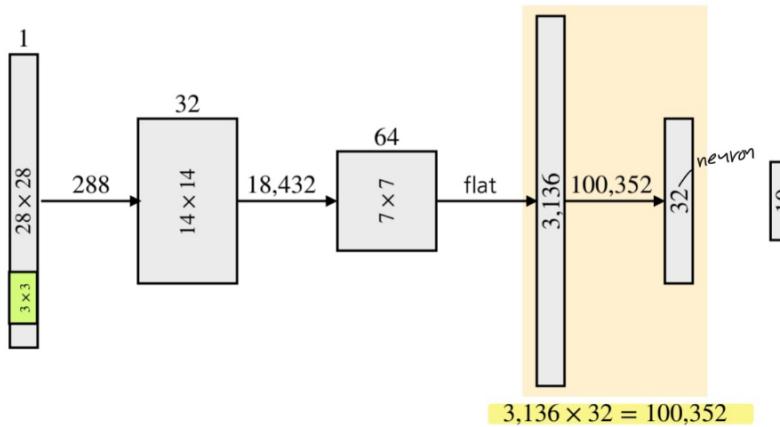
- Convolutions with 3×3 kernels, stride is two, and padding is one.
- Max-pooling with 2×2 kernels.

Parameter ≤ 288

20
Dept. of Science

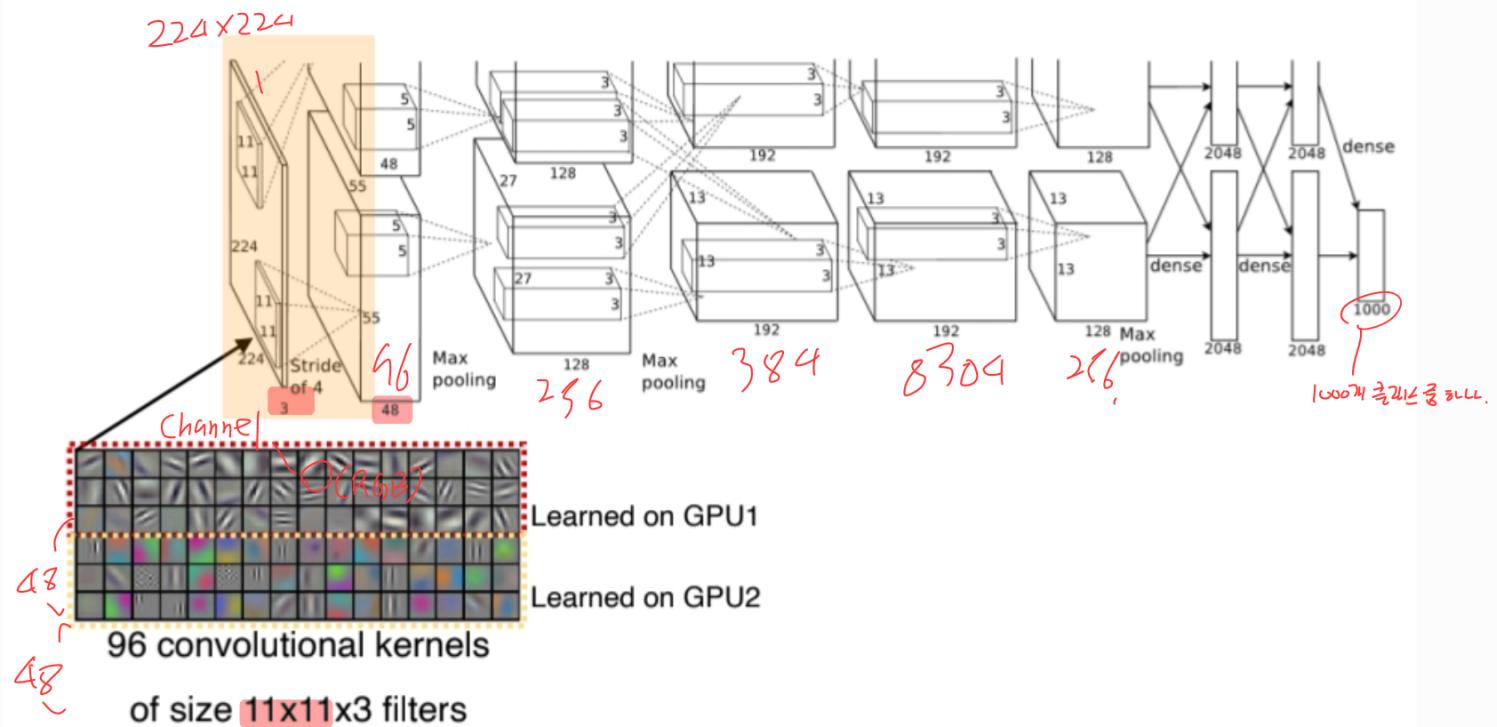


Convolutional Neural Network



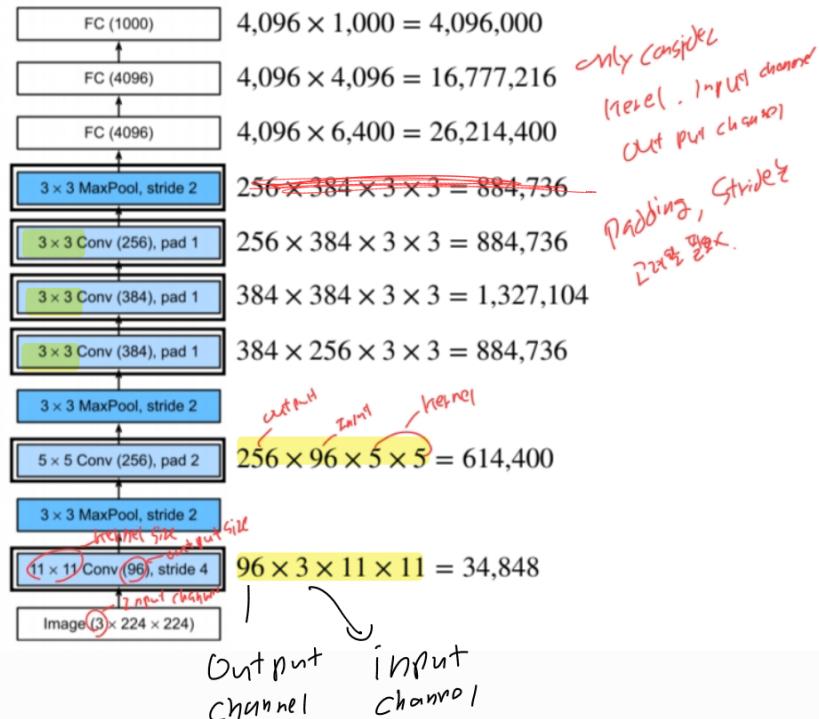
AlexNet

Architecture



Number of parameters

- Kernel: $[\text{out_channels} \times \text{in_channels} \times \text{kernel_height} \times \text{kernel_width}]$ 계산해보면 Can calculate number of parameters



Implementation Details

- Rectified Linear Unit (ReLU) activation
- Training on Multiple GPUs (2 GPUs) → spread network.
Put half of kernels on each GPU
- Local Response Normalization
- Overlapping Pooling
helps reduce overfitting by normalizing activation
- Reducing Overfitting
 - Data Augmentation
 - Dropout
Stride is less than pooling window size,
meaning the pooling windows overlap by one pixel,

VGGNet

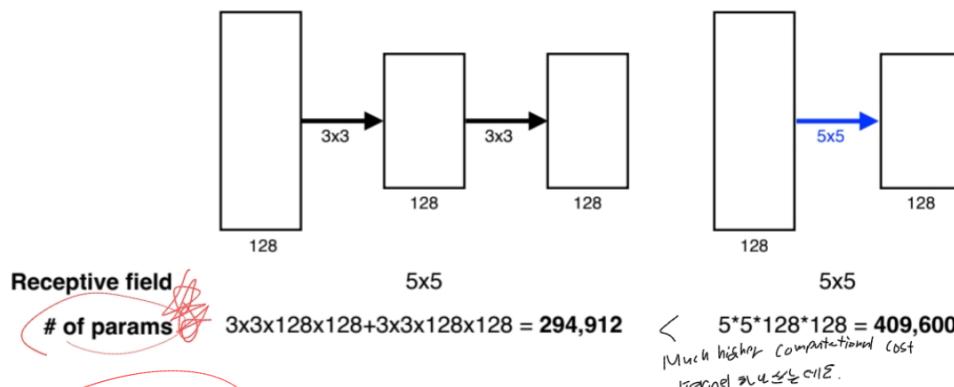
Enable deep networks. Using very small convolutional filters (3x3)

Basic building block.

- Convolutional layer with padding to maintain the resolution
- ReLU
- max-pooling to reduce the resolution

key idea is to use multiple convolutions

- Why use 3x3 convolution



GoogleNet

Multi-Branch Networks

stem: the first 2-3 convolution

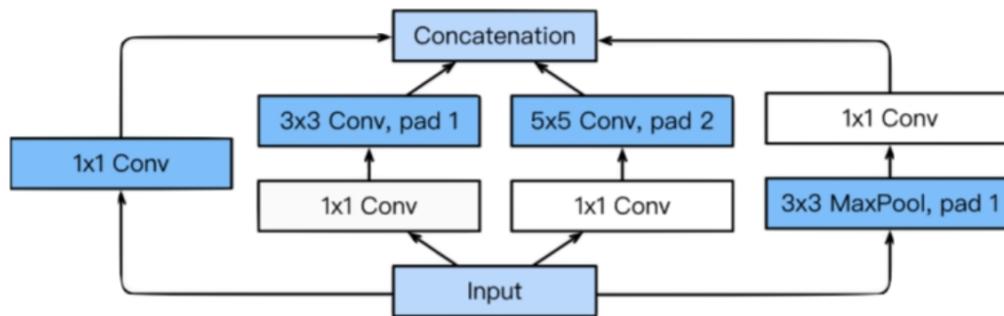
- body: convolutional blocks

head: maps the feature to the required problem

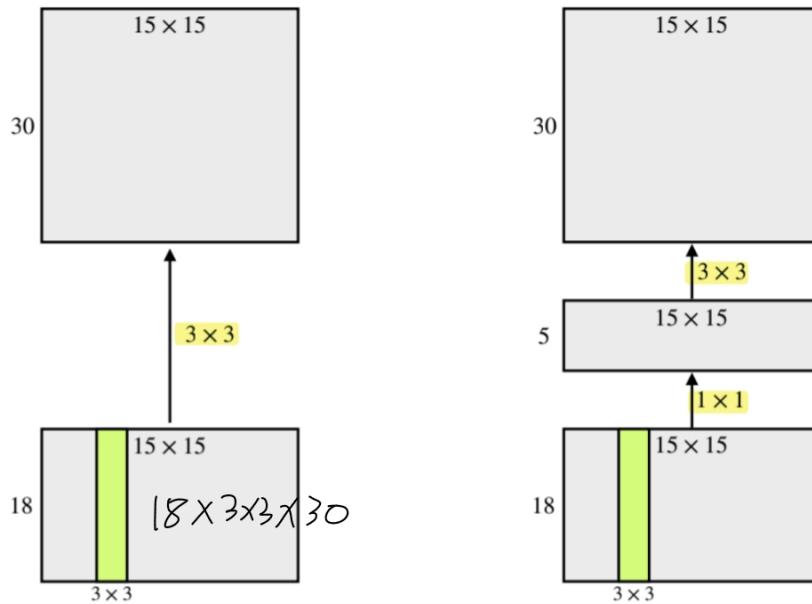
Inception blocks

reduce the number of parameters

1x1 Convolution can be seen as channel-wise dimension reduction

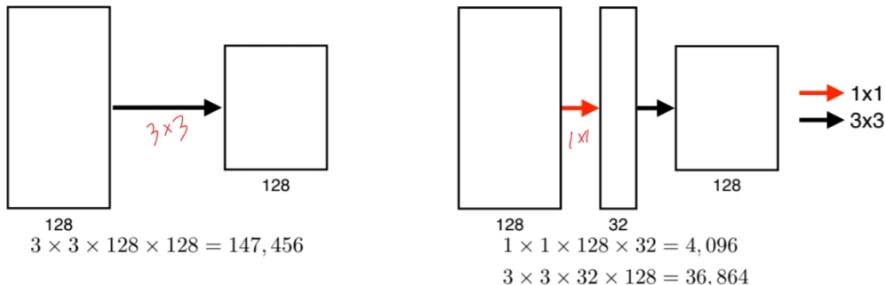


Why adding 1x1 Convolution?



$5 \times 3 \times 3 \times 30 =$

$18 \times 1 \times 1 \times 5 =$



1x1 convolution enables about 30% reduce of the number of parameters!

ResNet

- But, the **degradation** is not caused by overfitting, and **adding more layers** leads to higher training errors (and also higher test errors).

The key idea is that **every additional layer** should contain **identity function** as one of its elements.

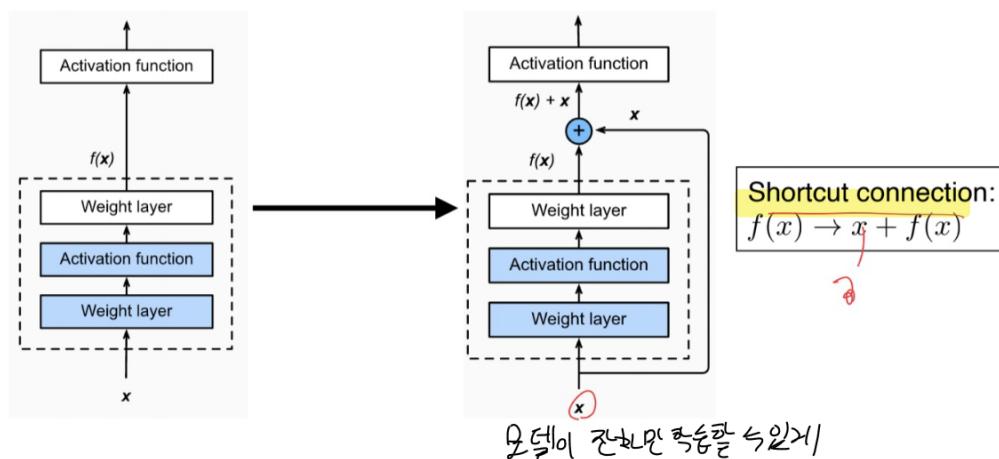
$$f(x) \rightarrow x, \text{ 원본을 그대로 두는 것}/$$

Residual mapping

Residual Blocks

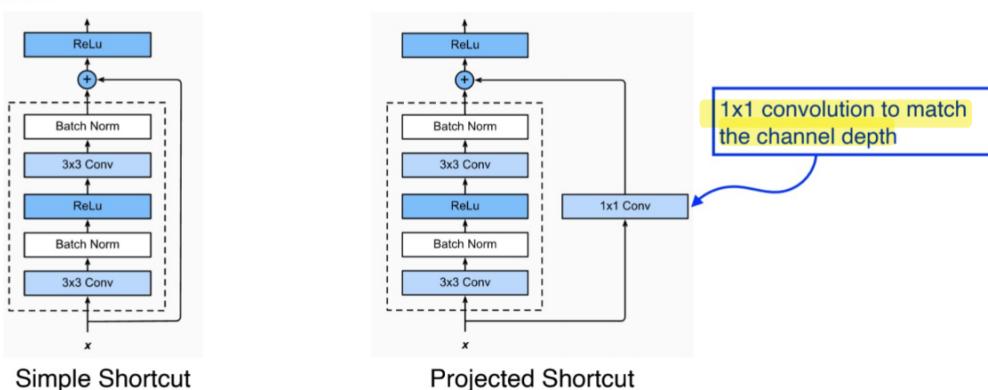


- Add an identity map (Shortcut connection)
- Shortcut allows the model to pass the input, and only learn the difference.



64

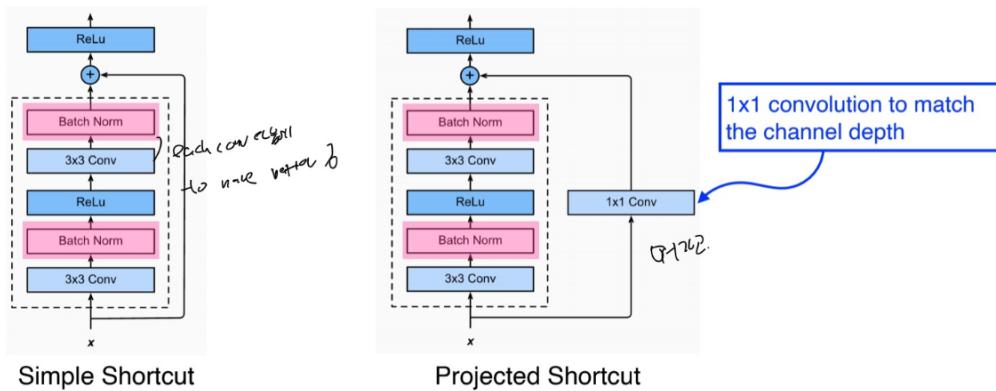
- Add an identity map after nonlinear activations.
- The projected shortcut uses a 1x1 convolution to match the number of channels.



65

Residual Blocks

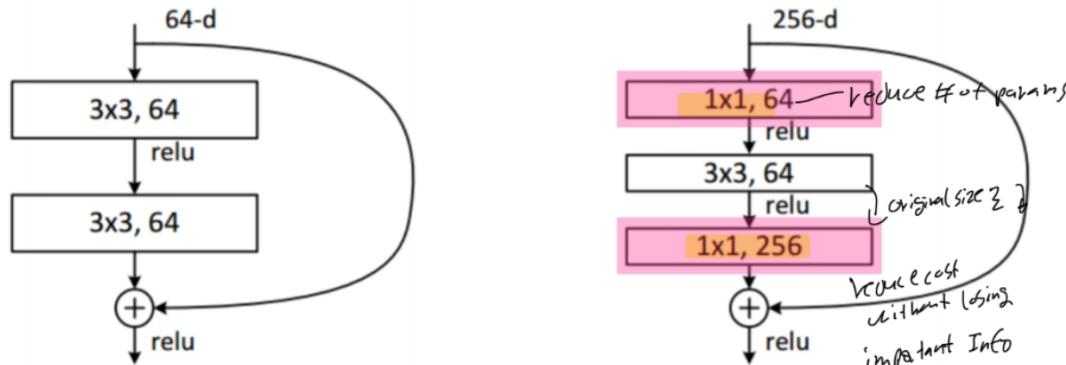
- Batch normalization after convolutions.



66

Bottleneck Blocks

- Similar to GoogLeNet's Inception architecture.
- $1x1$ convolution can be seen as channel-wise dimension reduction.



Performance increases while parameter size decreases.

Bottleneck Blocks

- Similar to GoogLeNet's Inception architecture.
- 1x1 convolution can be seen as channel-wise dimension reduction.

