

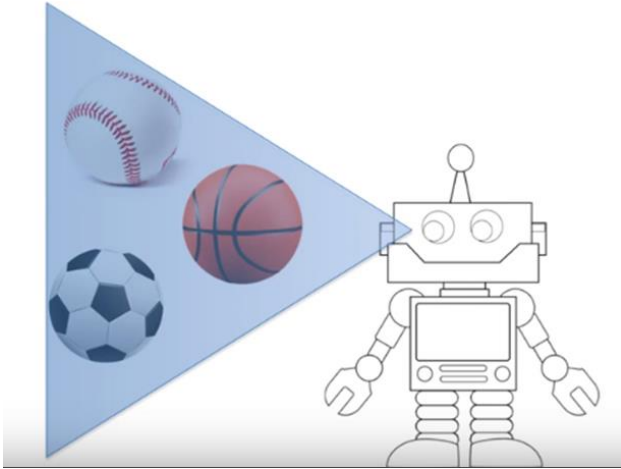


PNU Industrial Data Science

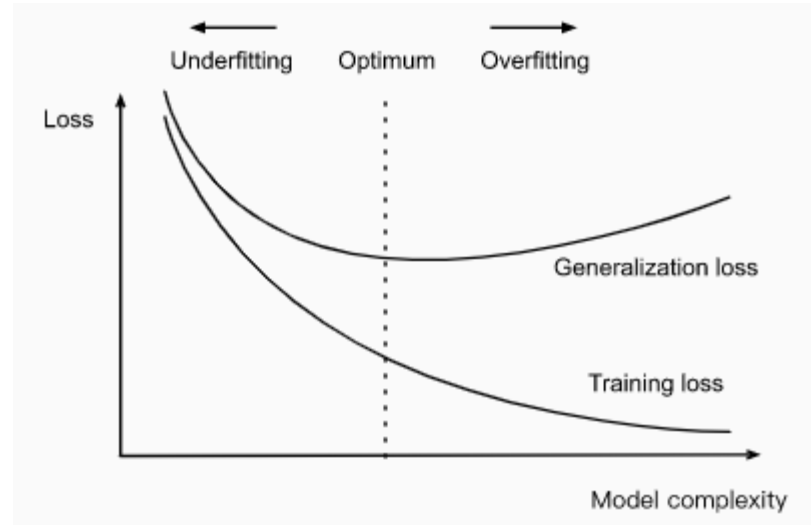
Convolutional Neural Network

합성곱 신경망

공을 구분하는 로봇



출처: <https://nittaku.tistory.com/289>



Contents

산업데이터과학은 산업현장에서 수집된 데이터를 분석하는데 필요한 기초 소양을 강의합니다.

01

Applications of NN

02

Deep Learning

03

CNN

Application of NN

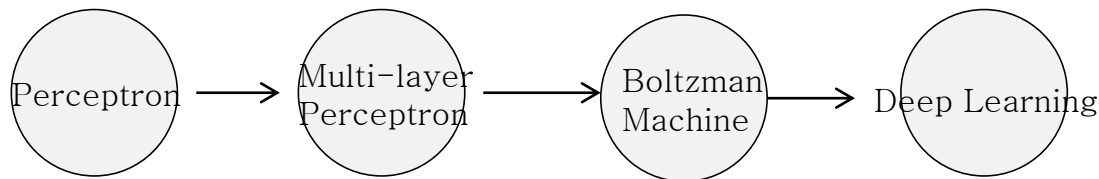
- Image recognition
 - Object recognition
 - Face recognition
- Voice recognition
- Spam-mail detection
- ...

Backgrounds of Deep Learning?

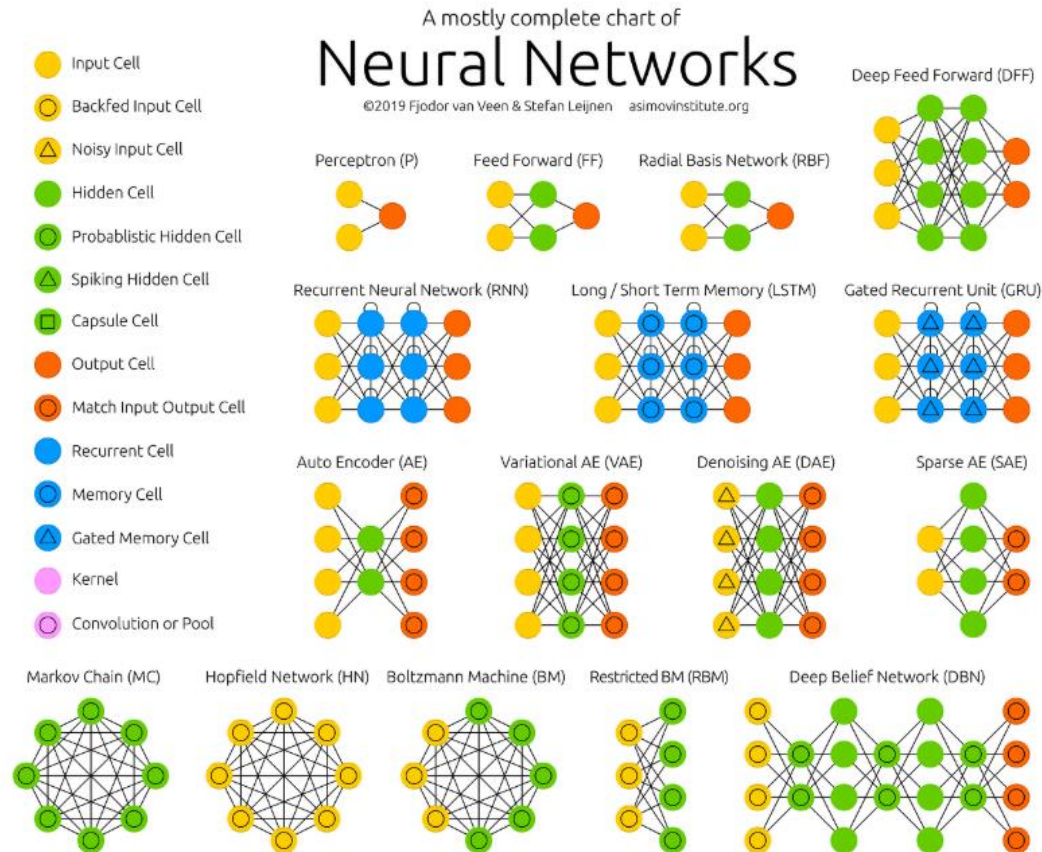
- Hard to program (?)
 - Automation (Decision system) vs. DSS
 - Problem is too complex
- Computing power
 - Big data
 - Mobile and SNS
- DL for business value
 - Ux
 - Advertisement
 - Recommendation

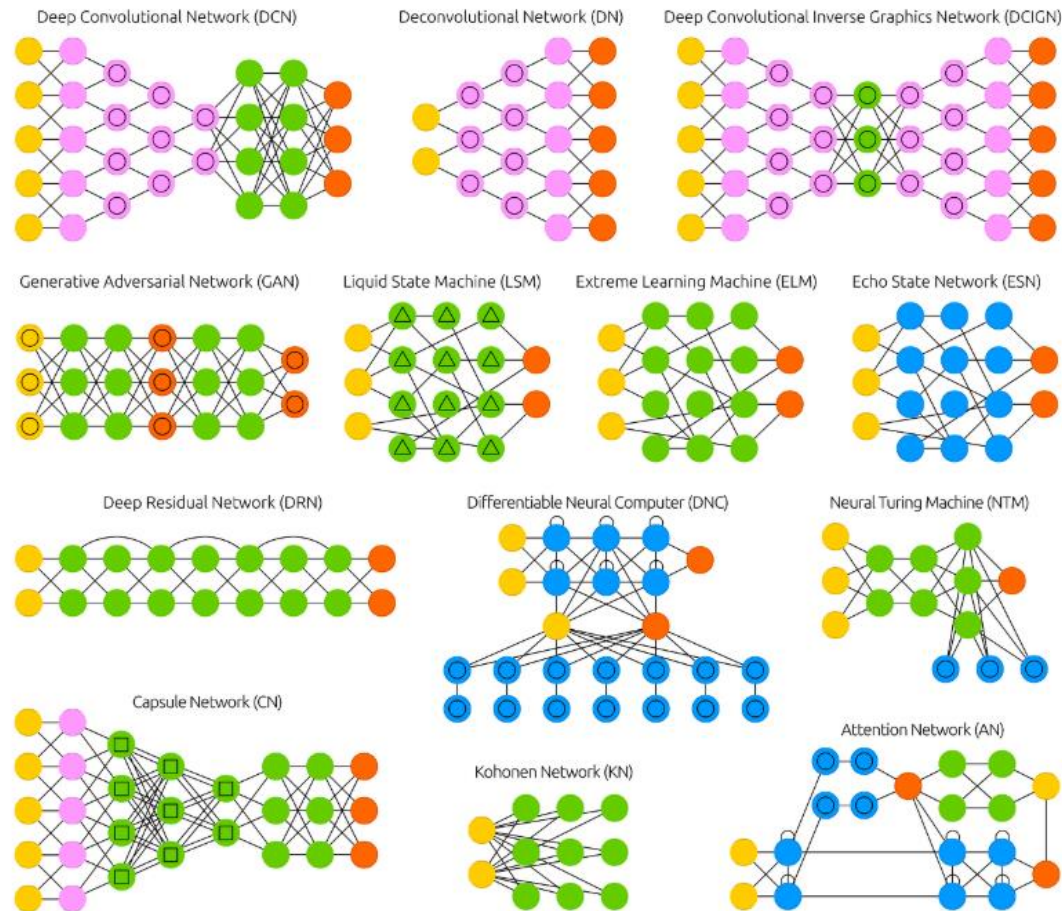
Brief history

- 1986 Back propagation algorithm
- 2004 DARPA Grand Challenge: Driverless car competition
- 2006 Deep Neural Networks by Hinton
- 2010 Turing award in learning theory (Valiant)
 - Probably approximately correct: Nature's Algorithms for learning and prospering in a Complex World
- 2011 Turing award in Bayesian networks (Pearl)
- 2012 DNNresearch (Hinton and Google)
- 2013 Quantum AI Lab for machine learning (Google)
- 2013 AI Labs for commercial ML research (Facebook)



Types of NN

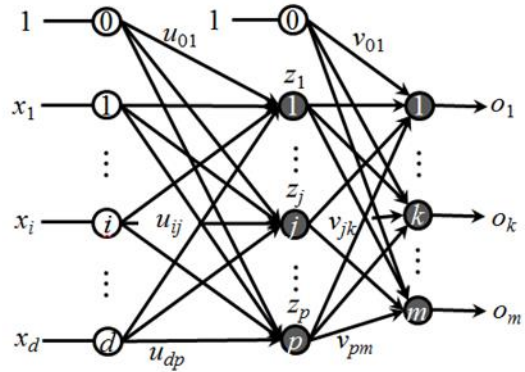




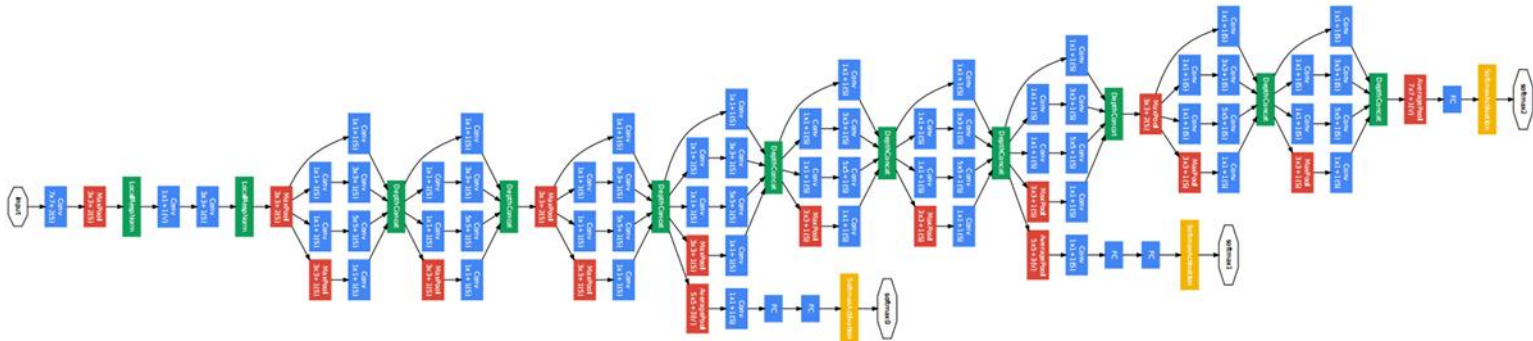
Source: <https://www.asimovinstitute.org/neural-network-zoo>

What is deep learning

- ANN



- DNN

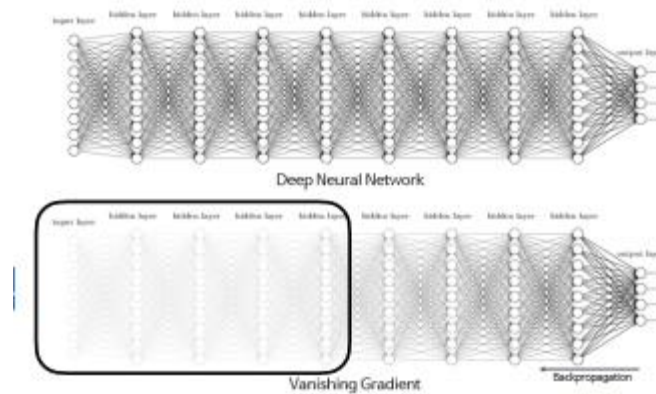
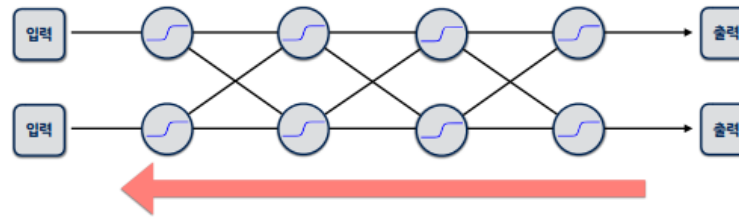


Three problems to be handled

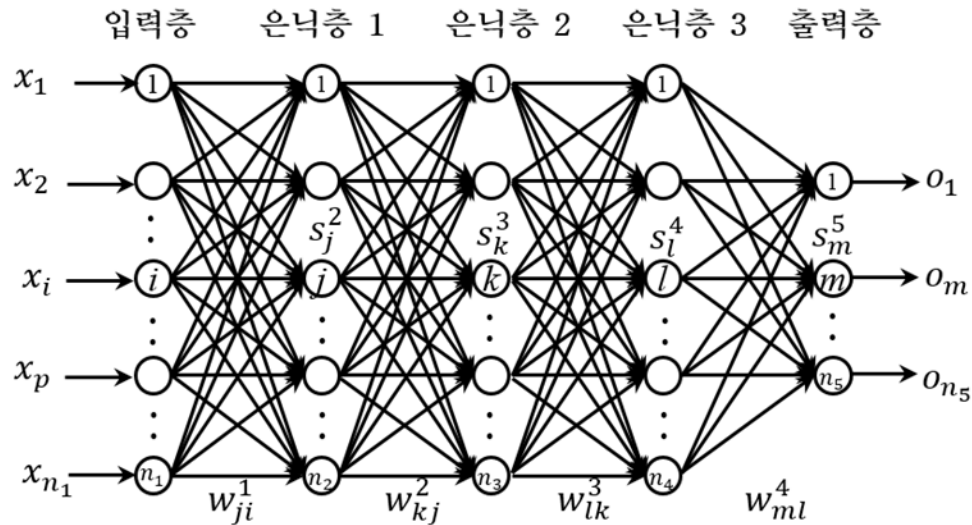
- Underfitting
- Overfitting
- Slow learning

1. Underfitting

- Learning method: Back propagation
 - Classifier produces too extreme result (if $f == s (= \sum_i w_i x_i)$)
 - Use Sigmoid function instead of step function

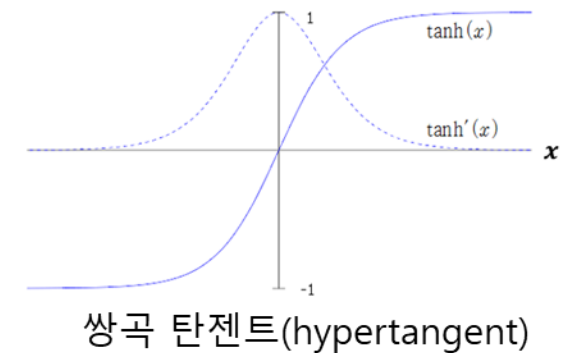
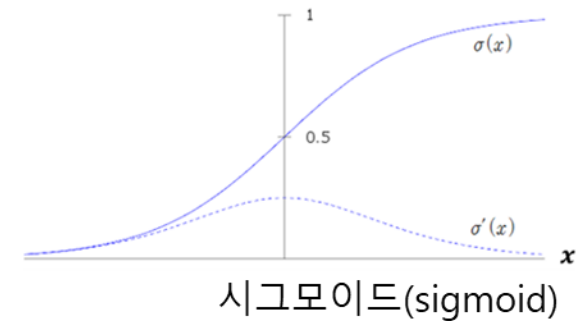


Problem of sigmoid BP



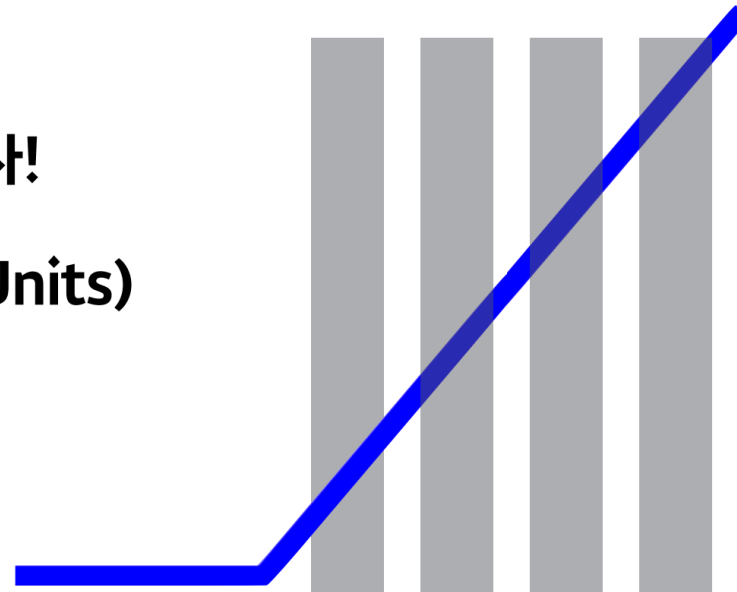
$$E = \frac{1}{2} \sum_{m=1}^{n_5} (y_m - o_m)^2$$

$$\frac{\partial E}{\partial w_{ji}^1} = \sum_{m=1}^{n_5} \sum_{l=1}^{n_4} \sum_{k=1}^{n_3} \sum_{j=1}^{n_2} (y_m - o_m) f'(s_m^5) w_{ml}^4 f'(s_l^4) w_{lk}^3 f'(s_k^3) w_{kj}^2 f'(s_j^2) x_i$$



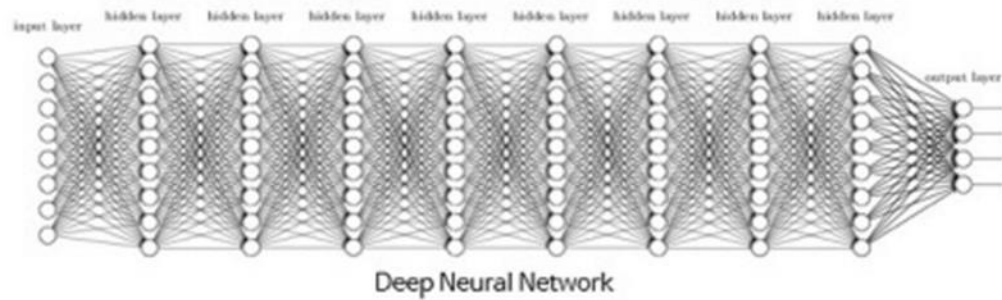
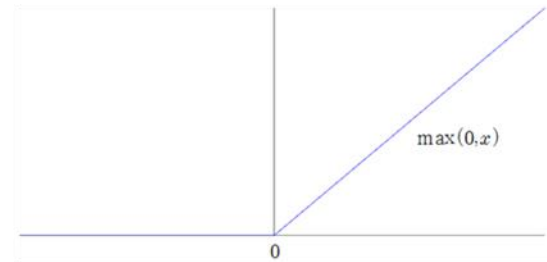
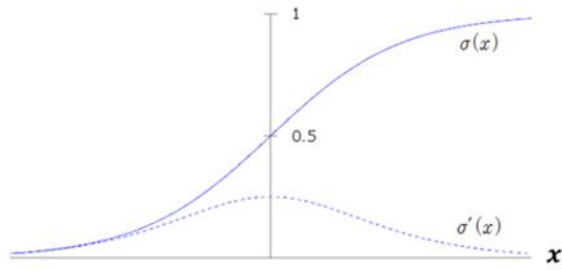
사그라드는 sigmoid대신
죽지않는 activation func을 쓰자!

→ **ReLU** (Rectified Linear Units)

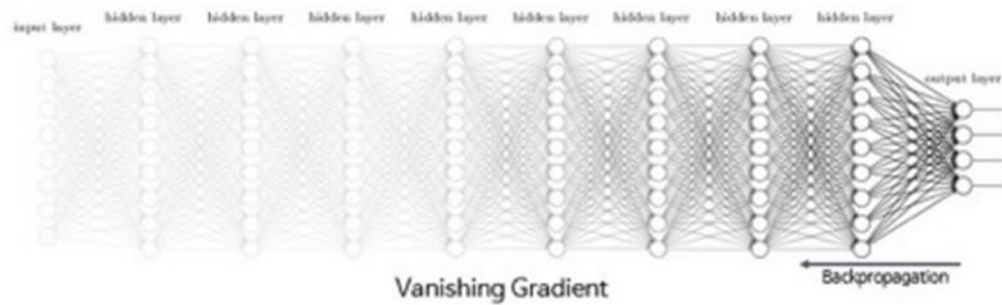


이녀석은 양의 구간에서 전부 미분 값(1)이 있다!

The effect of ReLU



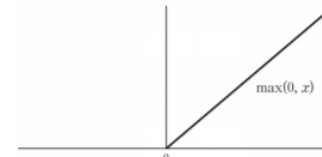
ReLU 사용



시그모이드 사용

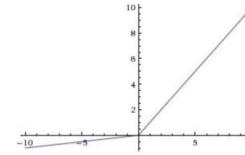
Variations of ReLU

- ReLU



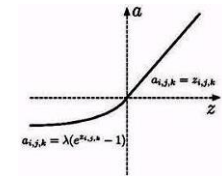
- Leaky ReLU

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



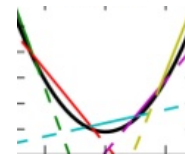
- ELU (exponential Linear Unit)

$$f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha(\exp(x) - 1) & \text{otherwise} \end{cases}$$



- Maxout

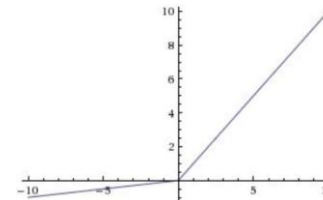
$$f(x) = \max_{i \in \{1, \dots, k\}} \{w_i^T x + b_i\}$$



- PReLU (parameteric ReLU)

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

α : Parameter learned



2. Slow learning

- Gradient Decent

$$\mathbf{W} \leftarrow \mathbf{W} + c(d - f)\mathbf{X}$$

$$\begin{array}{ccccccc} \text{Weigh update} & = & \begin{array}{c} \text{Direction} \\ \text{reducing err.} \\ \text{(descent)} \end{array} & \times & \begin{array}{c} \text{Size of one} \\ \text{step} \\ \text{(learning rate)} \end{array} & \times & \begin{array}{c} \text{slope} \\ \text{(gradient)} \end{array} \\ & & - & & \eta & & \nabla_{\theta} J(\theta) \\ & & & & & & \end{array}$$

- In big data
 - Weight update is too slow

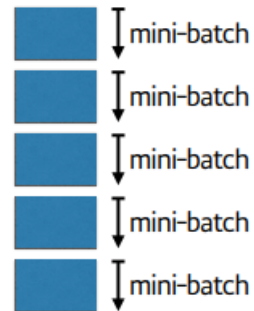
- SGD
 - Use mini-batch for updating the weight

Gradient Decent

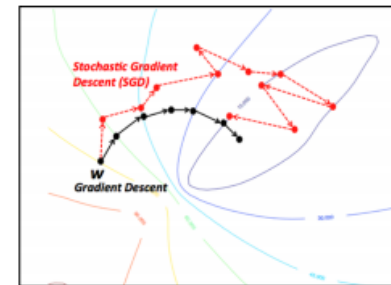


전부다 읽고 나서
최적의 1스텝 간다.

Stochastic Gradient Decent



작은 토막마다
일단 1스텝간다.



- Momentum vs. Adaptive

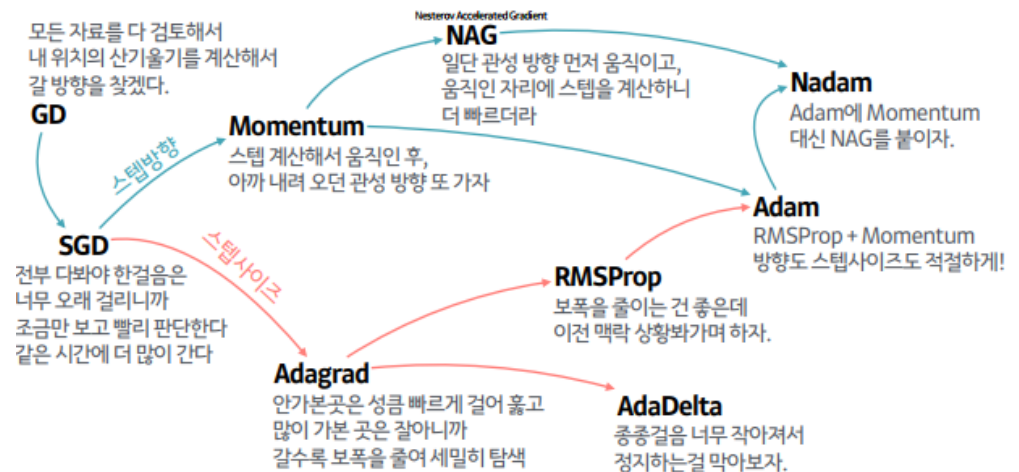
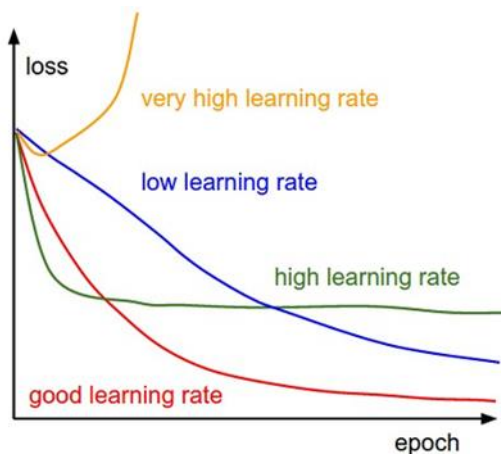
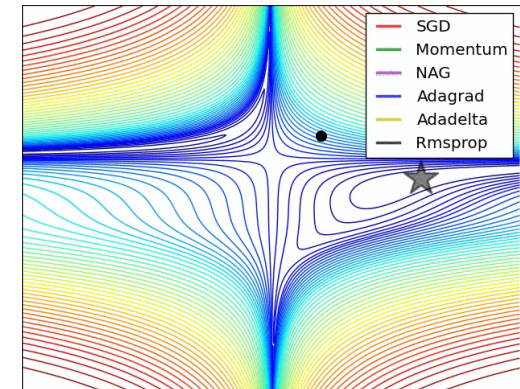
- Momentum $v_t = \gamma v_{t-1} + \eta \nabla_{\theta} J(\theta)$

$$\theta = \theta - v_t$$

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

- Adaptive

- Large step size for the variables with small change so far

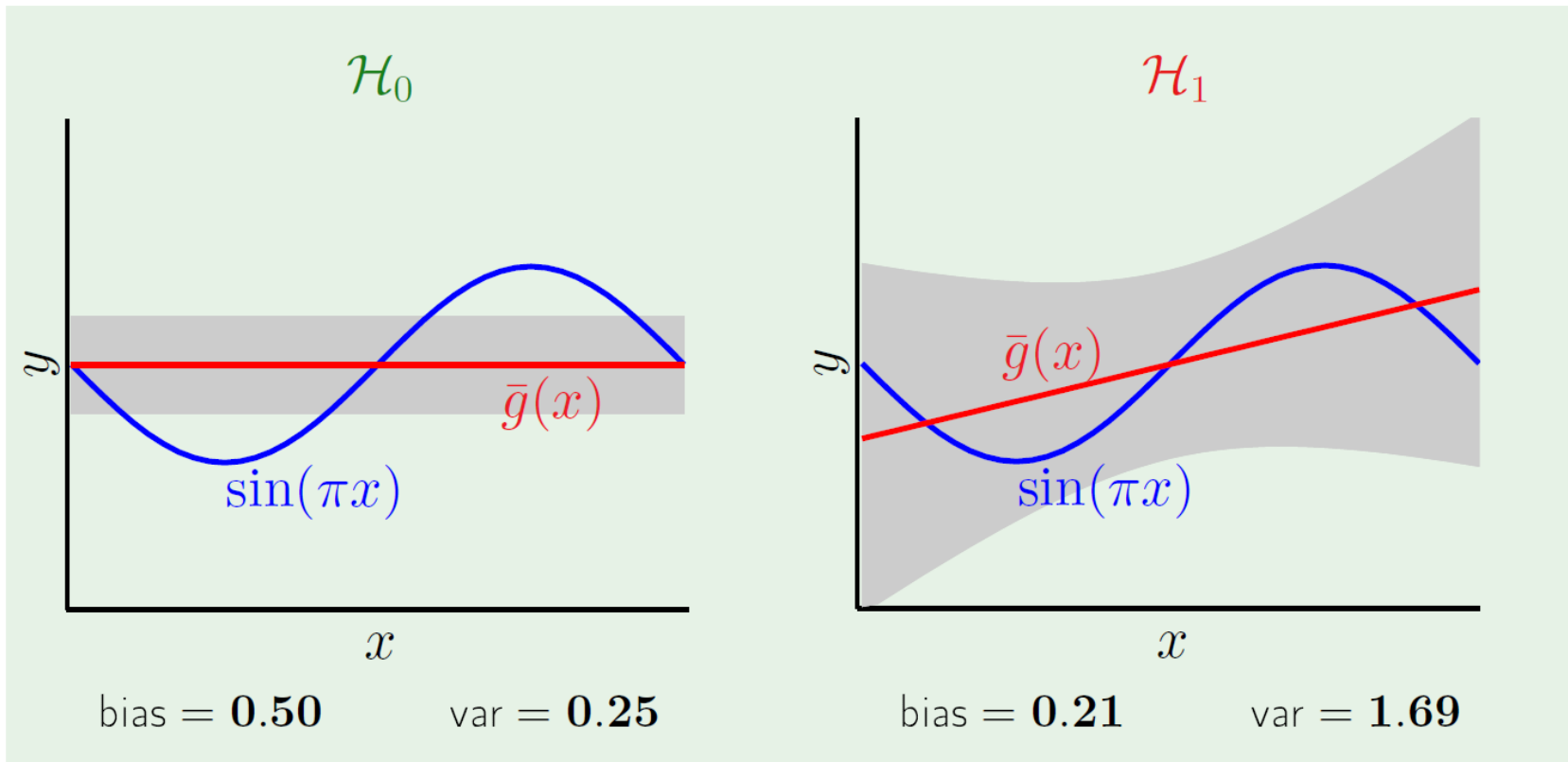


<https://seamless.tistory.com/38>

3. Overfitting

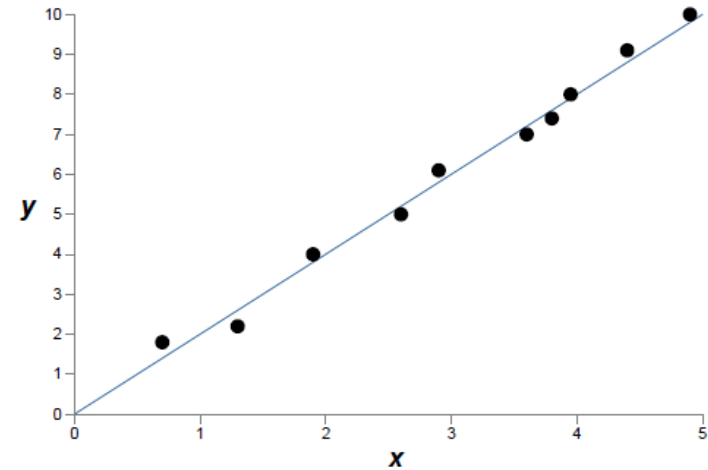
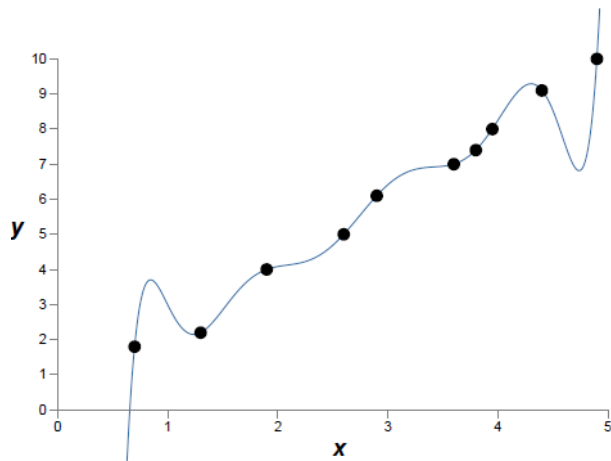
Bias-Variance Tradeoff

- Bias: Underfitting
- Variance: Overfitting



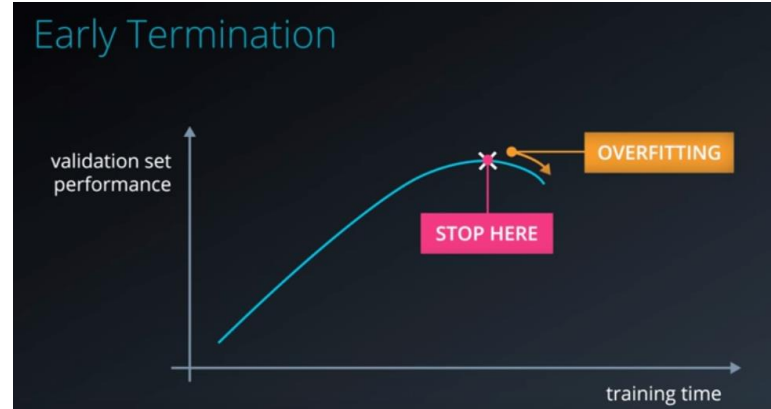
Overfitting

- Which one fits better?

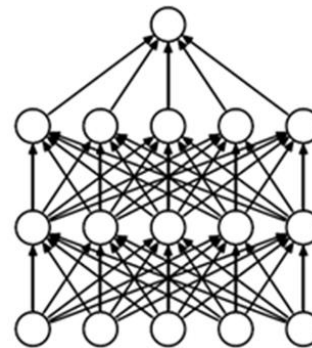


3. overfitting?

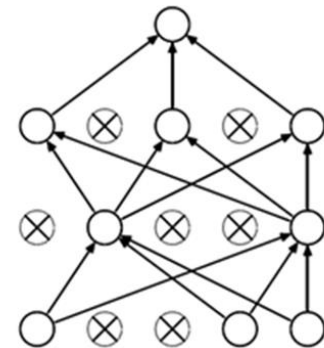
- Early Termination



- Regularization
 - L1, L2 normalization
 - Dropout
 - Remove some hidden nodes



(a) Standard Neural Net

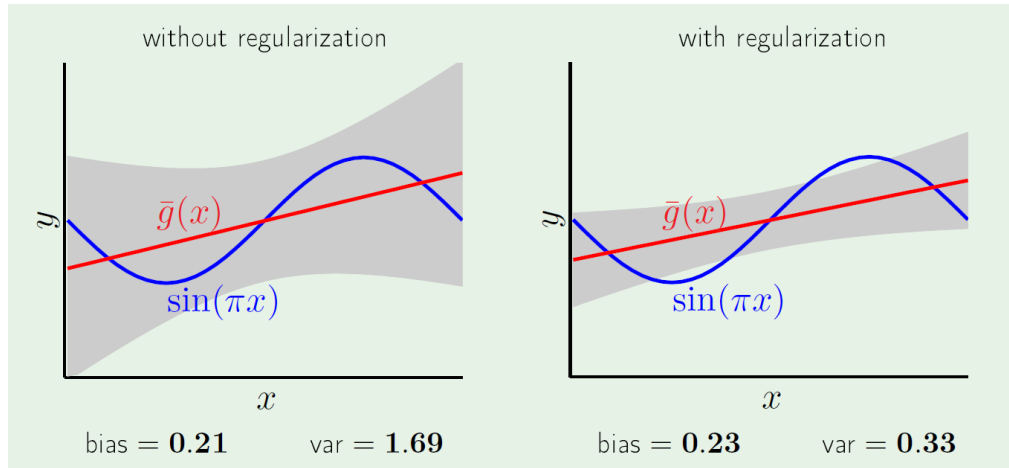


(b) After applying dropout.

Regularization

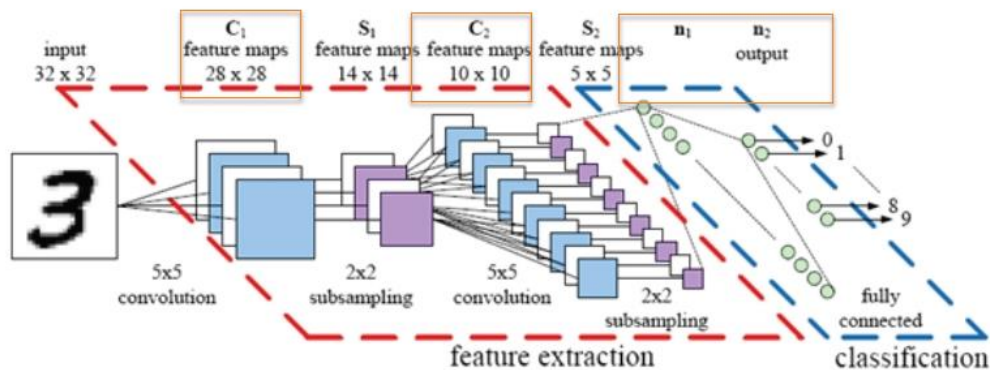
- L1 Regularization
 - L1 penalty uses the sum of the absolute weight.
 - $\min (Y - BX)^T(Y - BX) + \lambda\|B\|_1$
- L2 Regularization
 - L2 penalty is based on the squared weights
 - $\min (Y - BX)^T(Y - BX) + 0.5\lambda B^T B$

$$\begin{aligned}\text{Minimizing } E_{\text{aug}}(\mathbf{w}) &= E_{\text{in}}(\mathbf{w}) + \frac{\lambda}{N}\mathbf{w}^T\mathbf{w} \\ &= \frac{1}{N}(\mathbf{Z}\mathbf{w} - \mathbf{y})^T(\mathbf{Z}\mathbf{w} - \mathbf{y}) + \frac{\lambda}{N}\mathbf{w}^T\mathbf{w}\end{aligned}$$



Deep learning architecture

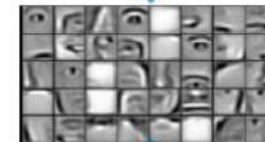
- CNN (Convolutional NN)
 - Repeat convolution and pooling



Feature representation



3rd layer
"Objects"



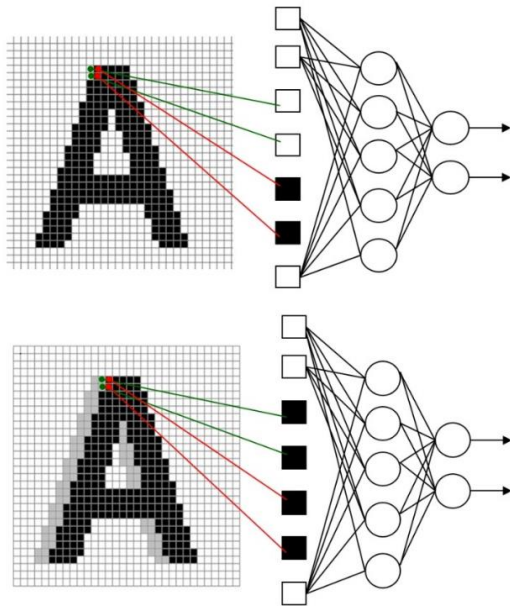
2nd layer
"Object parts"



1st layer
"Edges"



Pixels



Translation Invariance



Rotation/Viewpoint Invariance



Size Invariance



Illumination Invariance



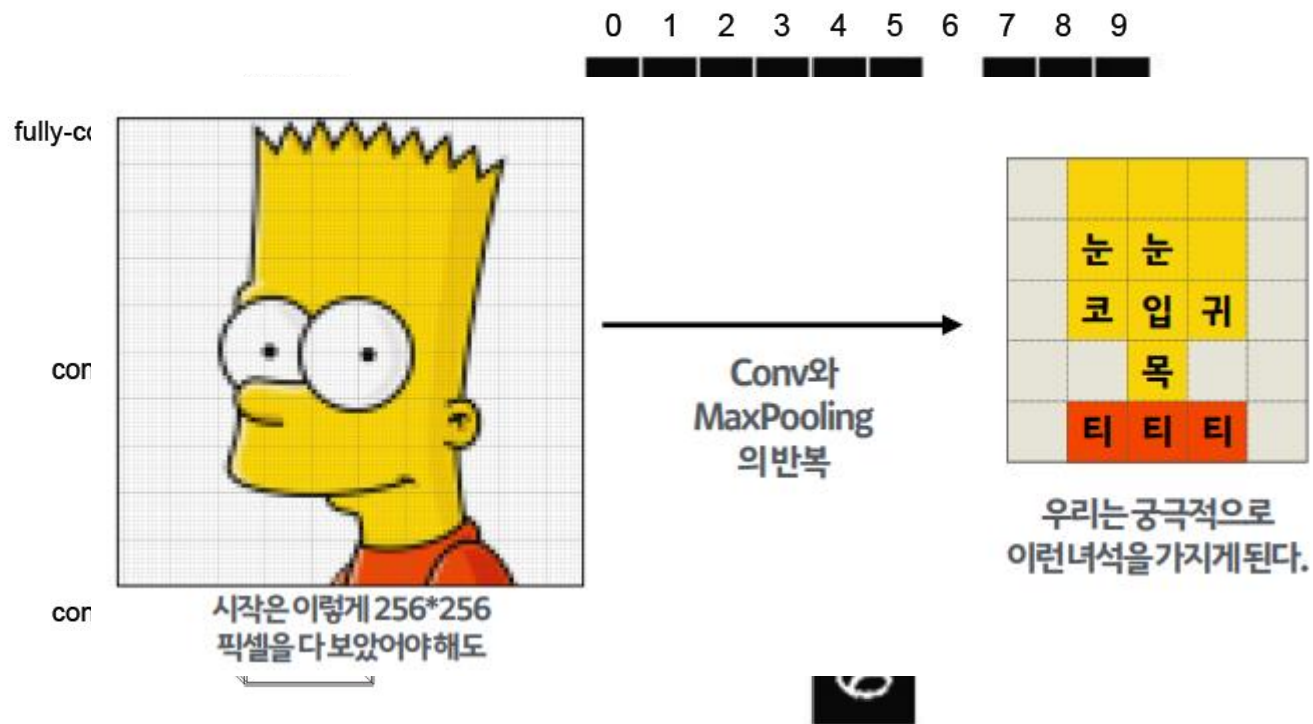
Matt Krause
mattkrause

- How do human know the classification?



Which one is cat and which is dog and how do you know ?

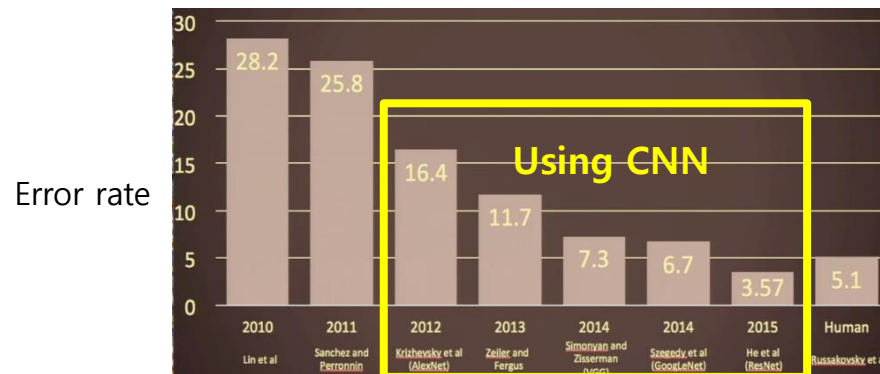
How computer can know ?



History of CNN

- CNN (Convolutional Neural Networks) was invented by Yann LeCun in 1998
- The first paper title: Gradient-Based Learning Applied to Document Recognition
- The popularity of CNN (Deep Learning) was extremely increasing after it always won in *ImageNet Challenge* (A computer vision world-cup) in 2012

For text recognition

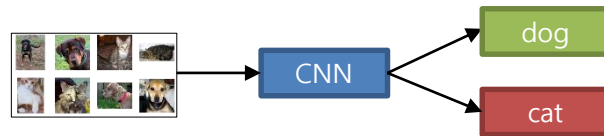


Fe-fei Li, et. al, CNNn23 Stanford Univ. 2017

Application of CNN

- Image Classification

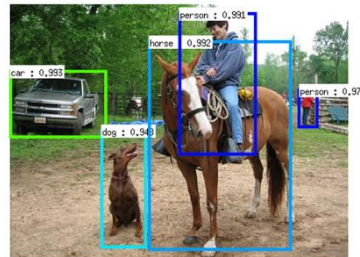
- OCR (Optical Character Recognition)



- Object Detection



- Image Generation



CNN Architecture

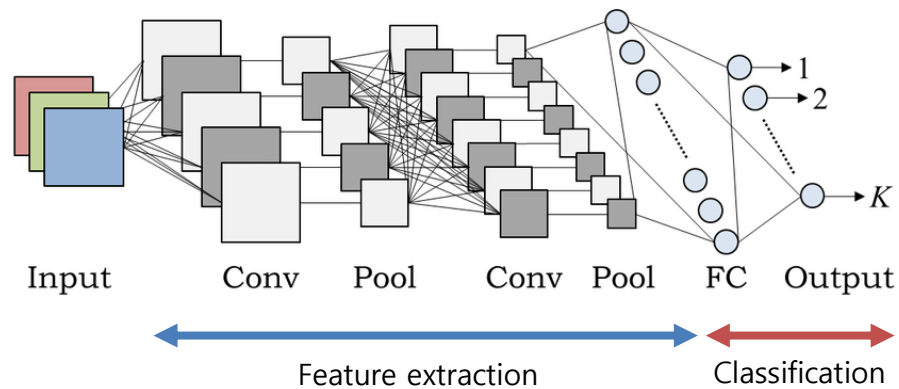
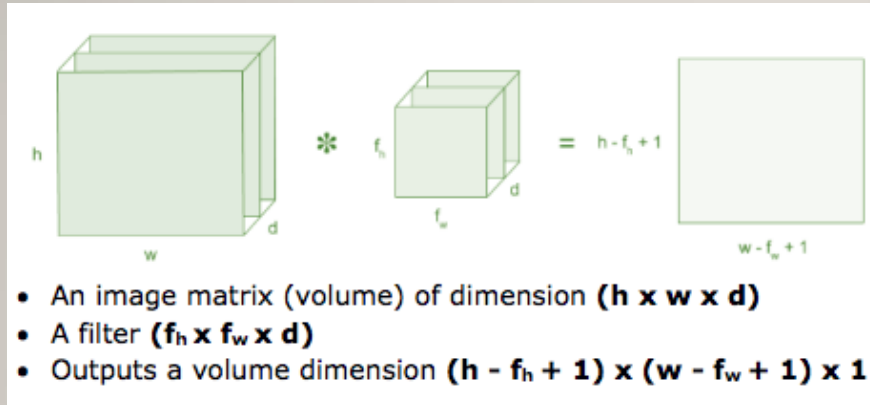
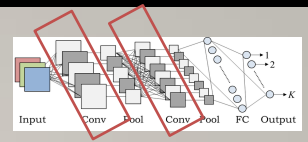
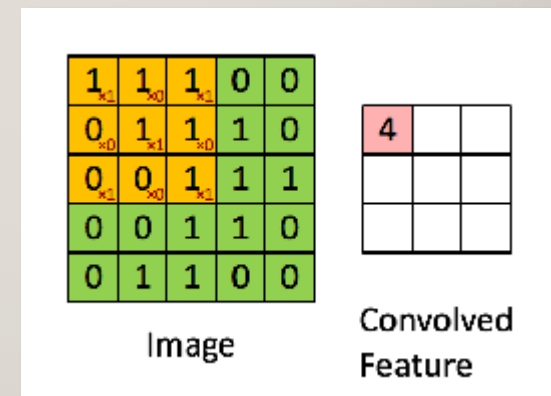
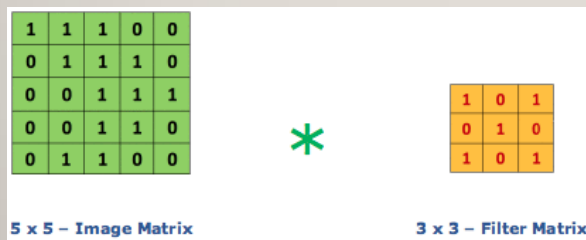


Image source: Akinori Hidaka, Takio Kurita, *Consecutive Dimensionality Reduction by Canonical Correlation Analysis for Visualization of Convolutional Neural Networks*

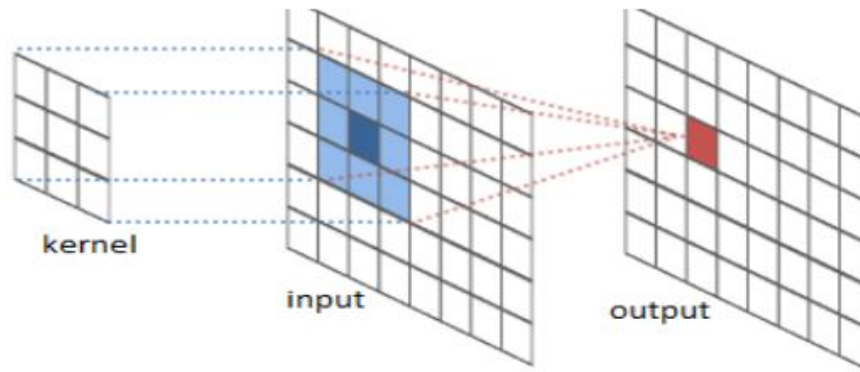
Convolution Layer (1)



- filter/kernel
- Receptive field



Convolution Layer (2)



Convolutional Layer (3)

The output of convolutional layer is the important features like edge, curve, dominant colour etc.

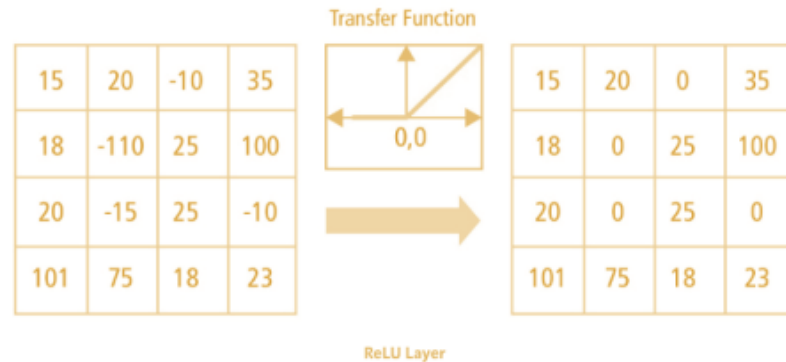
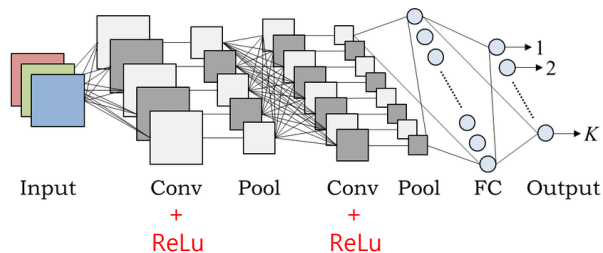
- Filter / kernels



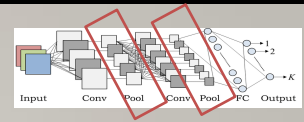
Operation	Filter	Convolved Image
Identity	$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$	
Edge detection	$\begin{bmatrix} 1 & 0 & -1 \\ 0 & 0 & 0 \\ -1 & 0 & 1 \end{bmatrix}$	
	$\begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$	
	$\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$	
Sharpen	$\begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$	
Box blur (normalized)	$\frac{1}{9} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$	
Gaussian blur (approximation)	$\frac{1}{16} \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$	

ReLU Activation Function

- ReLU stands for Rectified Linear Unit for a non-linear operation.
- $f(x) = \max(0, x)$



Pooling Layer



- Pooling layer applies non-linear downsampling on activation maps.
- Pooling is aggressive in discarding information

Max Pooling

29	15	28	184
0	100	70	38
12	12	7	2
12	12	45	6

2 x 2
pool size

100	184
12	45

Average Pooling

31	15	28	184
0	100	70	38
12	12	7	2
12	12	45	6

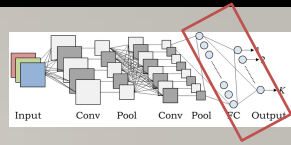
2 x 2
pool size

36	80
12	15

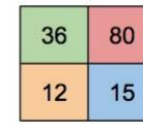
Types of pooling layer:

- Max Pooling
- Average Pooling
- Sum Pooling

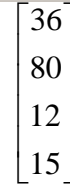
Fully Connected Layer (1)



- Neural Networks !
- **Flattened** the output from convolutional layer

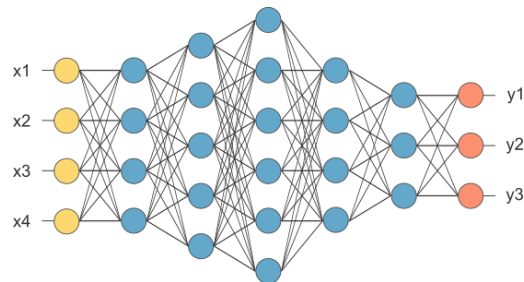


2D



Vector 1D

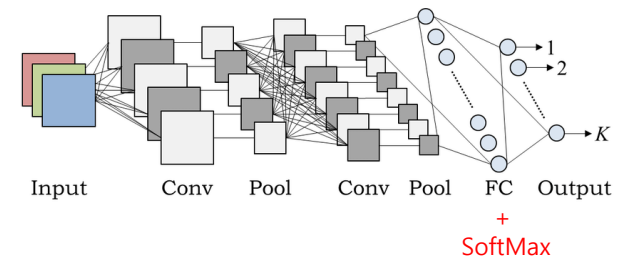
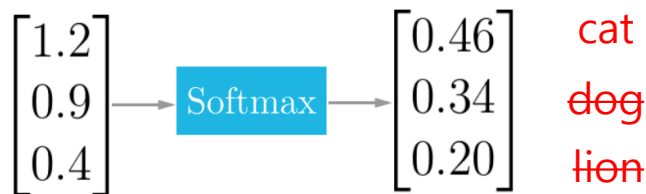
- The fully connected networks is used for feature classification
- The last layer uses **SoftMax** activation function (for classification)



How many neuron in each layer?

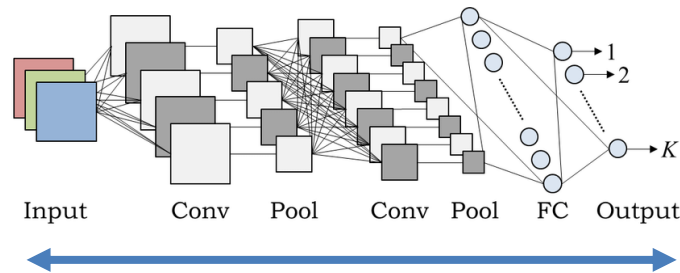
Fully Connected Layer (2)

- The softmax function squashes the outputs of each unit to be between 0 and 1
- The output of the softmax function is equivalent to a categorical probability distribution
- Index having highest value represent the output



Learning

- What dimension and kind of filter/kernel that we use ?
- What is the weight value that we use ?

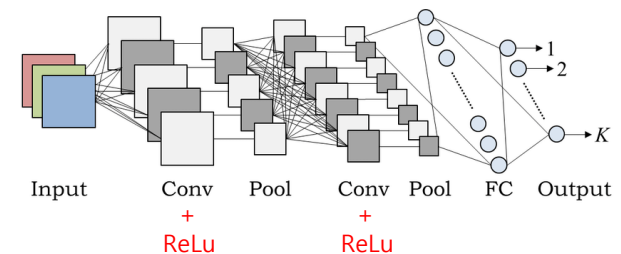


Back-propagation as supervised learning!

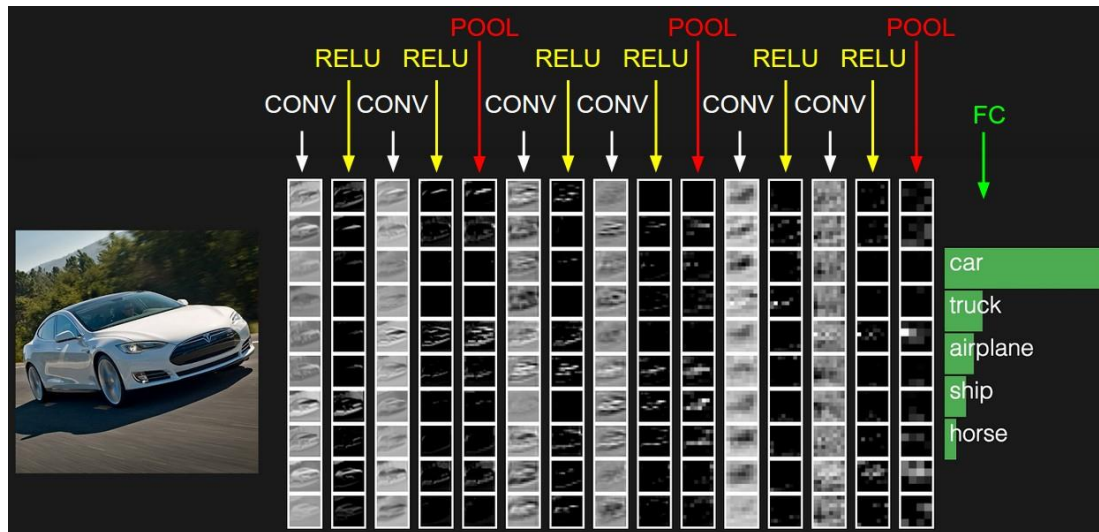
- We just need to initialize the value, and let the CNN through learning set the final filter and weight value itself.
- However, the number of layer and type of activation function are chosen based on your trial experiments.

Summary of How CNN works (1)

1. Provide input image into convolution layer
2. Choose parameters, apply filters with strides, padding if requires. Perform convolution on the image and apply ReLU activation to the matrix.
3. Perform pooling to reduce dimensionality size
4. Add as many convolutional layers until satisfied
5. Flatten the output and feed into a fully connected layer (FC Layer)
6. Output the class using an activation function (SoftMax) and classifies images.



Summary of How CNN works (2)



Advance of CNN

- How to choose the best architecture ? **it depends on our case**
- There is no standard architecture, like how many conv. Layer, pooling layer, strides, and number of fully connected layer. Do a trial experiment !