



# 1/18 Deep Feedforward Networks

## 1. Deep Feedforward Networks $\Rightarrow$ 여러 주요 application의 기본 ex) CNN, RNN, Transformer.

$\rightarrow$  가장 기본적인 neural nets

$\rightarrow$  예제: 하나의 벡터 (이미지, scalar)를 다른 벡터 mapping

$\rightarrow$  goal: approximate some function  $f^*$

$\rightarrow x \longrightarrow y$   $y$ 는 mapping이 뉴런치 기반으로 찾는 것!

ex) a classifier:  $y = f(x; \theta)$  maps input  $x$  to category  $y$

label Image

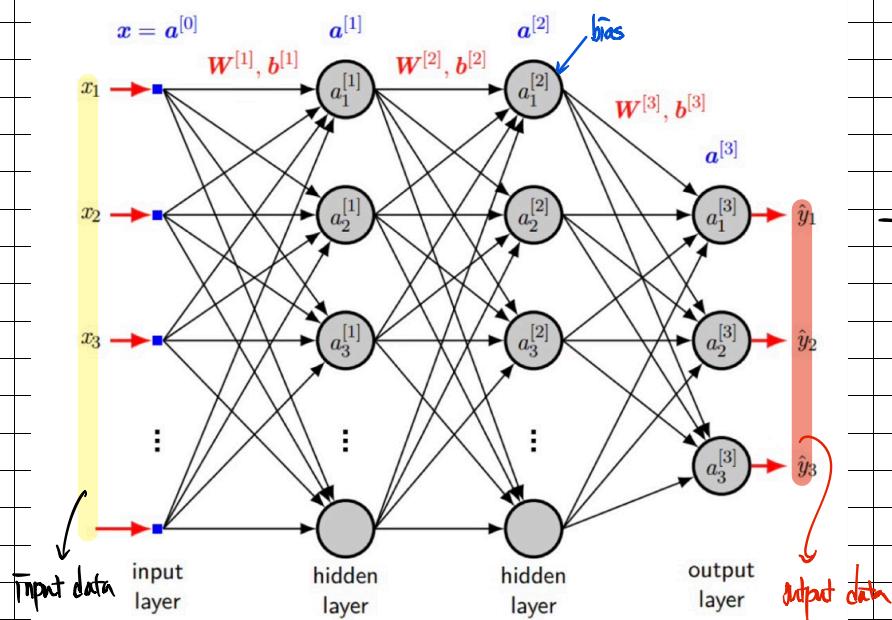
$\rightarrow$  how it works? : parameterize + learn

$\ni$  define a mapping  $y = f(x; \theta)$  and

$\ni$  learn parameter  $\theta$  that give the best approximation

## 2. Architecture

a feedforward neural net with two hidden layers.



$\therefore$  Input data (vector)  $\rightarrow f^*$   $\rightarrow$  output (vector)

$\rightarrow$  called feedforward because

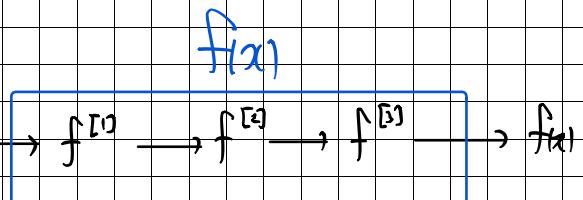
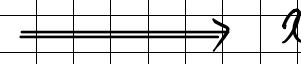
$\ni$  information flows  $x \rightarrow f \rightarrow y$

$\ni$  no feedback connection!  $\therefore$  ~~no~~  $\approx$   $x$

→ associated with a directed acyclic graph

⇒  $f^1, f^2, f^3$  connected in a chain!

$$f(x) = f^3(f^2(f^1(x)))$$



•  $f^{[l]}$ : called  $l$ -th layer

• final layer: called **output layer**

• chain length ⇒ model depth

3. Training neural nets : data를 이용해 우리가 원하는  $f^*(x)$  approximate 하는  $f(x)$ 를 구조는 고정!  
+ 첫 번째 계산하는 계정

• training: we drive  $f(x)$  to match  $f^*(x)$

• training data: noisy / approximate examples of  $f^*(x)$  ( $\epsilon$  term)

⇒ in supervised example  $(x, \underset{\text{label}}{y})$  with  $y \approx f^*(x)$

→ output layer 결과 개수 : MNIST, IMAGENET, BINARY

softmax

sigmoid

• behavior of the other layers: not directly specified by training data

⇒ these layers: called **hidden layers**

⇒ feature are distributed over hidden layers

↳ ?? : feature 한 층에 디자인으로 충돌!

여러 hidden layer가 distributed 하기 어렵다!

⇒ distributed representation

## 4. Understanding feedforward nets

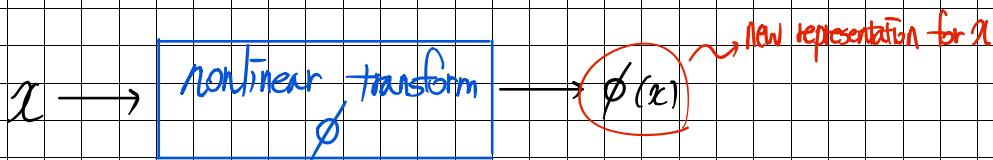
- begin with linear models

⇒ pros : efficient / reliable (closed form or convex)

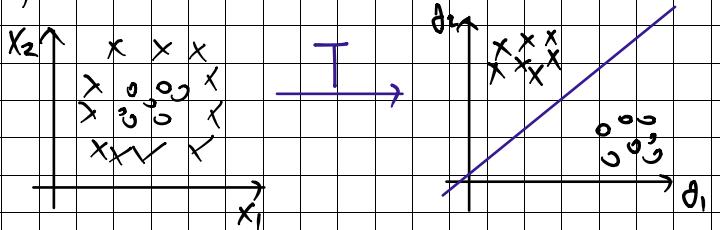
Cons : cannot understand interaction terms like  $x_1 \cdot x_2$

- to extend linear models to represent nonlinear function of  $x$

⇒ apply linear model transformed input  $\phi(x)$



Ex) SVM



- how to choose  $\phi$

In deep learning ⇒ Learn

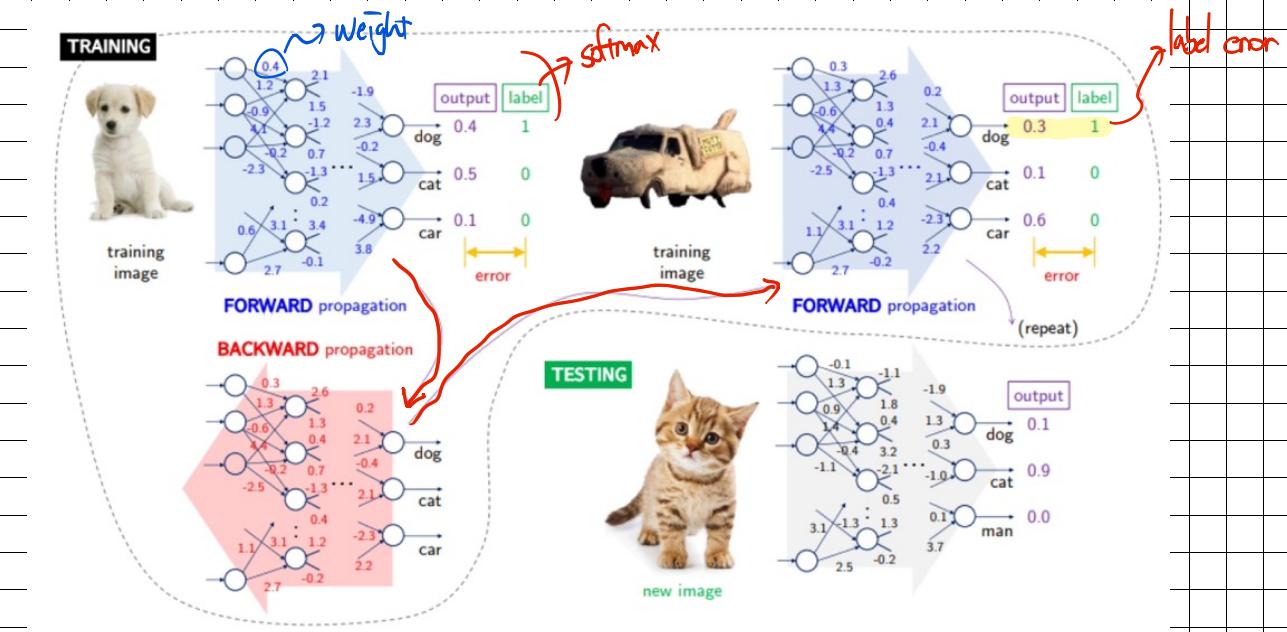
⇒ we have a model  $y = f(x; \theta, w) = \phi(x; \theta)^T w$

transform  $\phi$ 의 정의 문제  
: hidden layer #1 weight

$$x \xrightarrow{\phi(\cdot; \theta)} \text{feature} \xrightarrow{w} y$$

$\theta$   $y$   
Output mapping's parameter : output layer weight

## 5. Concepts of training & testing a neural net



## 6. Modern Neural Nets.

→ the same backprop / gradient descent : still use ...

biggest change ?!

loss function  $\frac{\partial L}{\partial w}$  or  $\frac{\partial L}{\partial b}$  ?!  
big  $\frac{\partial L}{\partial w}$  or  $\frac{\partial L}{\partial b}$  ?!

1. MSE → cross entropy loss

2. Sigmoid → ReLU

gradient vanishing

## 7. Cost function for Neural nets.

- total cost function

$\rightarrow$  primary cost function + regularization term

$$\Rightarrow \text{min } E_{\text{avg}} = \min \left( E_{\text{train}} + \frac{\lambda}{n} \|w\|_2^2 \right)$$

regularization term

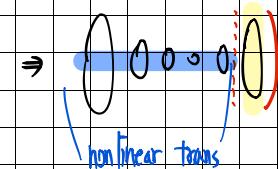
- most modern neural nets : trained using maximum likelihood

i.e. cost function =  $-\text{NLL}(\hat{y}|y)$

= cross entropy ( $\hat{y}|y$ )

- a recurring theme : gradient of cost function must be large  $\approx 1$  / predictable ?

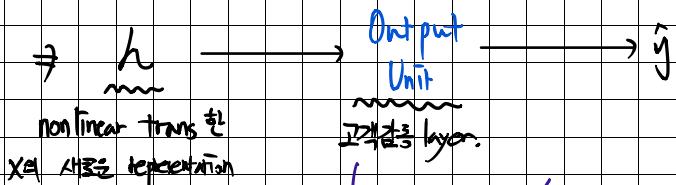
## 8. Output units



transform of  $\phi$   $\Rightarrow$   $\hat{y}$   
hidden layer  $\neq$  weight

- Suppose : a feedforward net provides hidden features  $h = f(x; \theta)$

- Output layer : provides "additional" transformation from  $h$

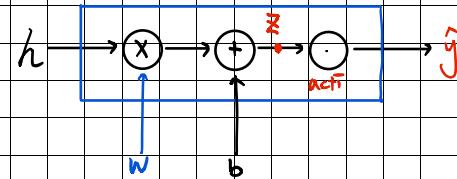


use linear / sigmoid / softmax  
binary multi

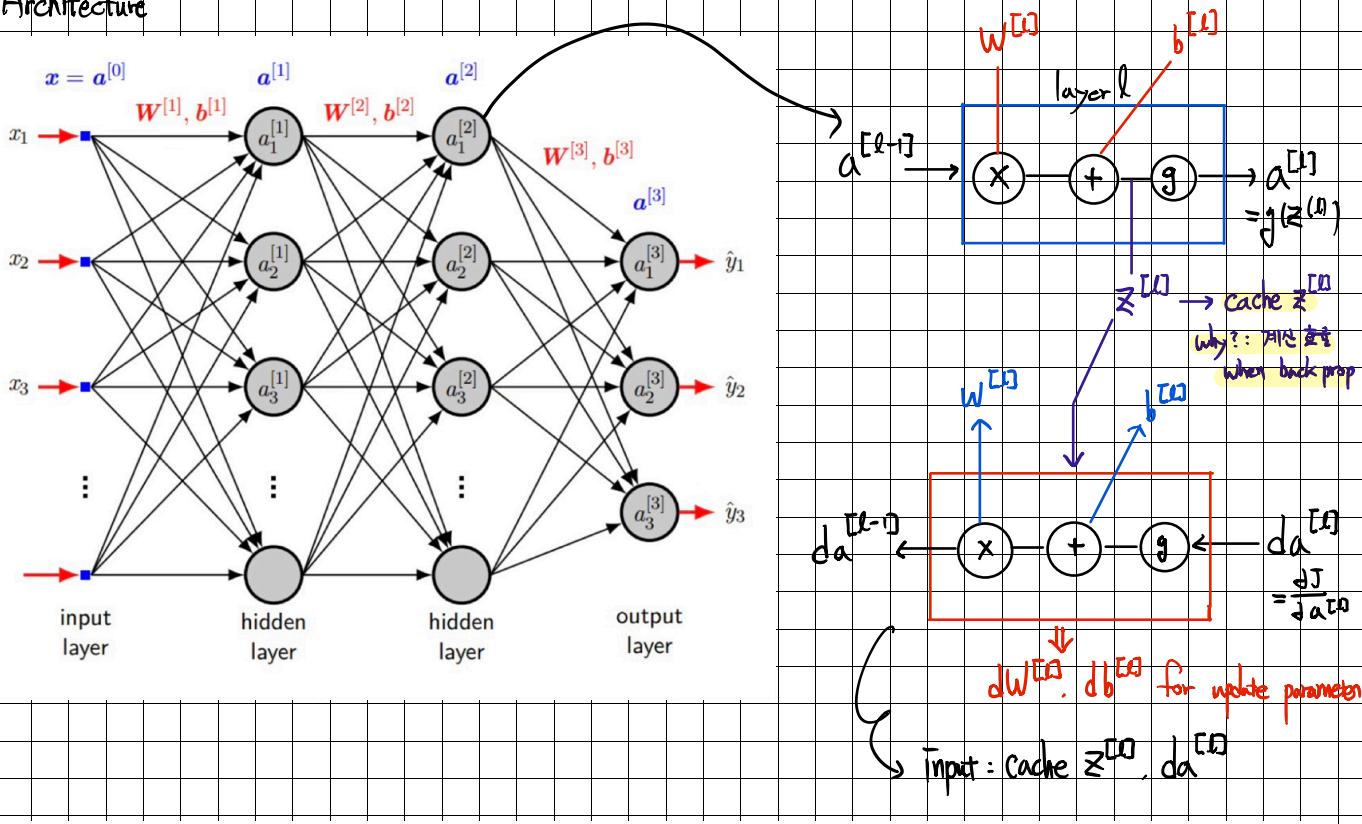
distribution	# classes	# trials (samples)
Bernoulli	2	1
multinoulli	$k$	1
binomial	2	$n$
multinomial	$k$	

## 9. Type of Output Units

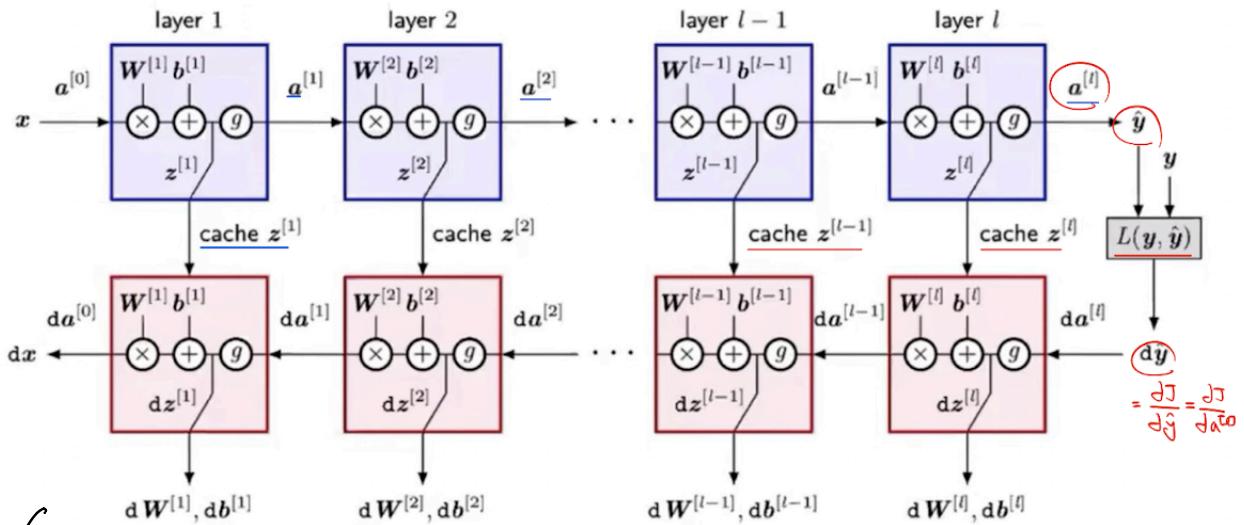
type	output	formula	output distribution
GLM or A (Link function.)	linear	$\hat{y} = W^T h + b$	Gaussian
	sigmoid	$\hat{y} = \sigma(w^T h + b)$	Bernoulli
	softmax	$\hat{y} = \text{softmax}(W^T h + b)$	multinoulli
ex) exp	$\hat{y} = \exp(W^T h + b)$	Poisson	



## 10. Architecture



## Overall architecture

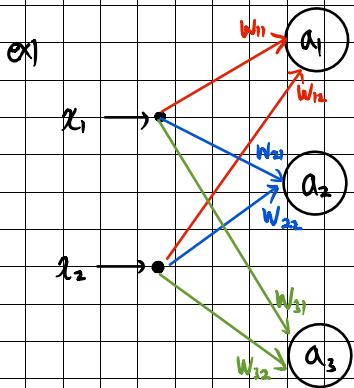


by gradient we get  $\Rightarrow$  parameter update!

$$W^{[l]} \leftarrow W^{[l]} - \epsilon dW^{[l]}$$

$$b^{[l]} \leftarrow b^{[l]} - \epsilon db^{[l]}$$

## 11. Weight matrix conventions



• BL (Right-Left) convention

$$\overleftarrow{W} = \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \\ w_{31} & w_{32} \end{bmatrix} \xrightarrow[3 \times 2 \text{ } 2 \times 1]{} \text{then, } \overleftarrow{W} \vec{x} \longrightarrow \vec{z} = \overleftarrow{W} \vec{x} + b$$

• LR convention

$$\overrightarrow{W} = \begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{bmatrix} \xrightarrow[3 \times 2 \text{ } 2 \times 1]{} \text{then, } \overrightarrow{W}^T \vec{x} \longrightarrow \vec{z} = \overrightarrow{W}^T \vec{x} + b$$

## 12. Hidden Units

- What they do ?!

1. Accept a vector of input  $\vec{x}$

2. Compute affine transformation  $\vec{z} = \vec{W}^T \vec{x} + b$

$\Rightarrow$  linear transformation  
+  
translation

3. apply an element-wise nonlinear function  $g \xrightarrow{\text{nonlinear}} \vec{z}$

$$g(\vec{z}) = \begin{bmatrix} g(z_1) \\ g(z_2) \\ g(z_3) \end{bmatrix} \text{ when } \vec{z}^T = [z_1, z_2, z_3]$$

4. return activation  $\vec{z} = g(\vec{z})$

$\xrightarrow{\text{e.g.}} \text{ReLU (rectified linear units)}$

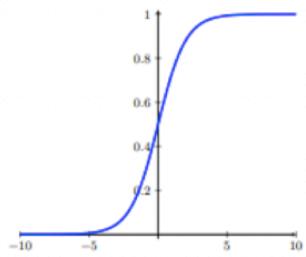
• hidden units differ only in activation function

$\Rightarrow g$  isn't an affine trans. or etc.

사실 다른 nonlinear function 사용 가능.  
ex: sin, cos...  
but published only if clearly significant ↑

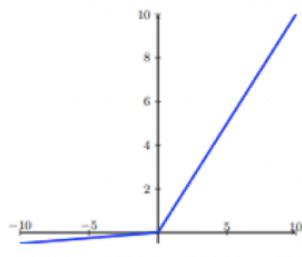
• mixing activation function types in a layer is uncommon

### 13. Activation functions



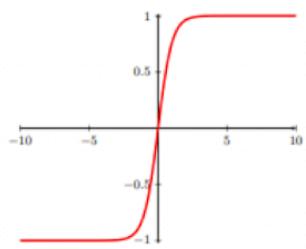
• sigmoid

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



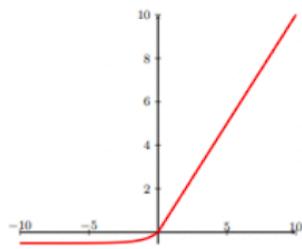
• Leaky ReLU

$$\max\{0.01z, z\}$$



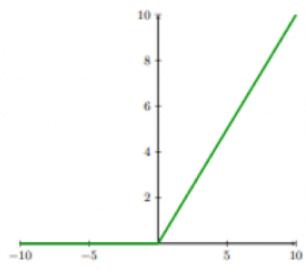
• tanh

$$\tanh(z)$$



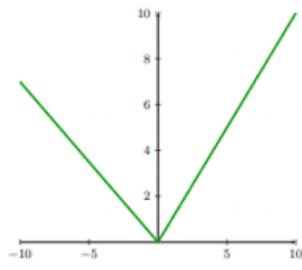
• ELU

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



• ReLU

$$\max\{0, z\}$$



• maxout

$$\max\{z_1, z_2\}$$

$$\vec{z} = \vec{W}^\top \vec{x} + b \quad . \quad \text{ReLU}(\vec{z}) = \max\{0, \vec{z}\} \quad = g(\vec{z})$$

$$g'(\vec{z}) = \begin{cases} 1, & \vec{z} \geq 0 \rightarrow (+)-\text{region: grad} = 1 \\ 0, & \vec{z} < 0 \rightarrow (-)-\text{region: grad} = 0 \end{cases}$$

• Pros !!

i) no saturation in (+) region

ii) converges faster than sigmoid

: (+) region with gradient > 1 편평하지 않음

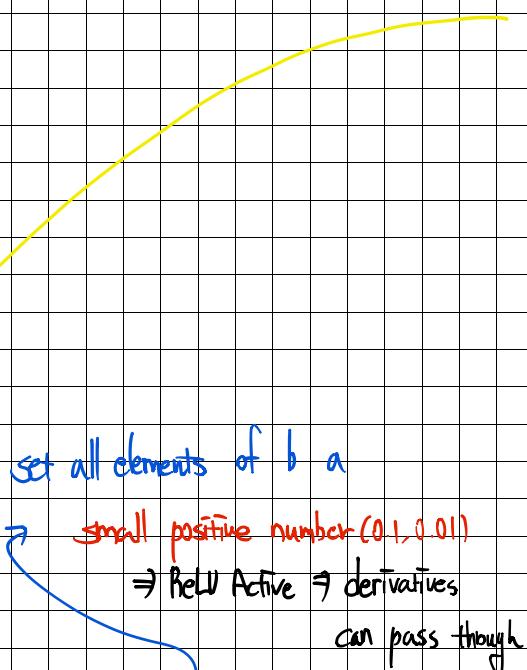
• Cons !!

i) not zero-centered output

→ 편평한 양수로 차우침...

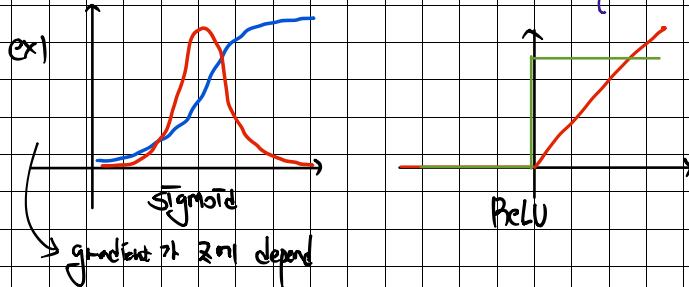
→ gradient가 편평함으로 누적

→ optimization이 zig-zag (둘리 철렁)



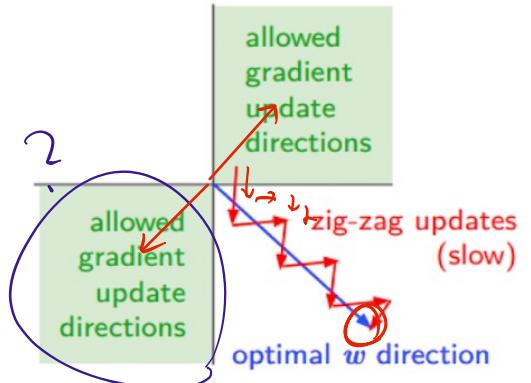
## 4. ReLU optimization : why ReLU?

- easy to optimize ( $\rightarrow$  so similar to linear units)
- only when ReLU Active
- derivatives through ReLU
  - (1) remain large when ReLU Active
    - $\Rightarrow$  Active  $\Leftrightarrow$  gradient = 1  $\forall z$ .
    - $\Rightarrow$  Sigmoid / tanh  $\Leftrightarrow |z| \uparrow \rightarrow$  gradient = 0  $\forall z$ .
    - $\Rightarrow ?!$  Vanishing gradient  $\forall z$ !
  - (2) not only large but consistent
    - $\Rightarrow$  second derivative: 0  $\forall z$  at training of  $\forall z$ ?



## 5. One-sided inputs slows down training

- Precall :  $\frac{d}{dw} (wx) = x$
- $\Rightarrow$  weight update direction  $\in \mathbb{R}$   $\Leftrightarrow$  not zero-centered
- One-sided input??
  - all (+), (-),  $Ex \neq 0$
  - $\Rightarrow$  not zero-centered



• Then why one-sided inputs are matter?

$\Rightarrow$  slow convergence!  $\Rightarrow$  zig-zag updates

$\Rightarrow$  update direction biased

So ! Normalization  $\xrightarrow{\text{ZG}}$  ! (Batch Norm?)  
 Input ( $\pm 1$  year)

## 16. Exploiting Linearity : gradient at ~~the~~<sup>the</sup> ~~end~~<sup>end</sup> off!

- principle of ReLU

⇒ models are easier to optimize if their behavior is close to linear  
⇒ non-linear ⇒ gradient vanishing problem  
flow of gradient

- cons...: linear boundary



## 17. Architecture exploration

- main architecture considerations in chain-based architectures.

⇒ network depth and layer width

more depth

more row? node

# of layer

- a feedforward net with a single layer

⇒ sufficient to represent any function

⇒ may fail to learn generalize correctly

node ↑ → estimate weight & bias ↑

→ # parameter ↑

INN!

- deeper network ⇒ increase depth!

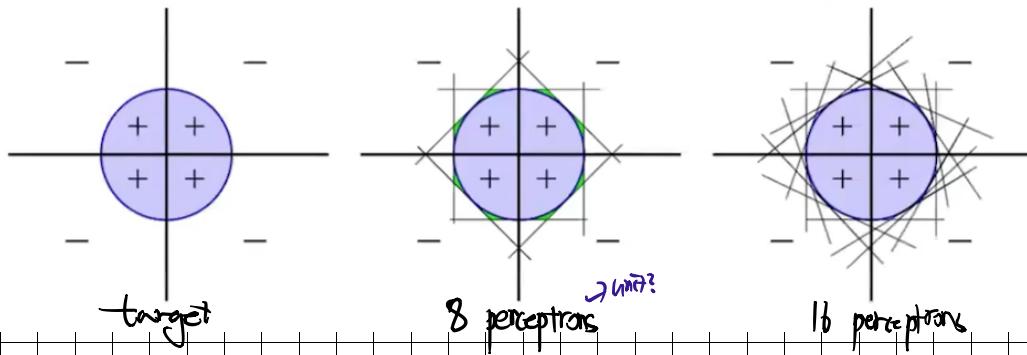
⇒ use far fewer units per layer and fewer parameters.

⇒ generalize better to test set

⇒ but harder to optimize (vanishing, exploding gradient)

## 18. Universal approximation Theorem

- a feedforward net with linear output layer + hidden layers  
⇒ can approximate any function (give enough **hidden units**)



## 19. Network size

- UAT says there exists a network large enough to achieve any accuracy  
⇒ but does not say how large this network will be.
- Unfortunately an **exponential number of hidden units** may be required in the worse case
- empirically: greater depth  $\Rightarrow$  better **generalization**