# CS 461: Machine Learning Principles

Class 16: Oct. 28

Back Propagation Algorithm &

Its Application to Deep Neural Net (DNN)

Instructor: Diana Kim

$$y = \underline{w}^t \cdot \phi(x)$$

# Linear Modeling & <u>Finding Global Minimum</u>

So far, we have studied the functions for the task of regression and classification. They were all linear functions $\dot{y} = w^t x$.
And algorithms to learn the functions from data.
(MMSE, Logistic MLE, and SVM dual)

- Linear Regression: $y = w^t x$

$$W* = \arg\min_w J(W) = \arg\min_w ||y - \Phi w||^2$$
$$= \arg\min_w (y^t - w^t \Phi^t)(y - \Phi w)$$
$$= \arg\min_w ||y||^2 - 2y^t \Phi w + w^t \Phi^t \Phi w$$

$6 \cdot y = w^t \cdot x$

$y = w_k^t \cdot x$

$f^{g \cdot L}$

$\sim \iota o o (\tilde{y})$

$J(w)$ based data

$\cdot \nabla^2 \cdot J(w) \geq 0$

$\cdot \nabla J(\overset{*}{w}) = 0$

$\Phi^t \cdot \Phi \cdot w = \Phi^t \cdot y$

# Hessian of Objective Function (Loss) & Convexity

- $\nabla^2 J(w) \geq 0 \leftrightarrow J(w)$ is convex. $\leftrightarrow J(w)$ has a global minimum.

- if $\nabla^2 J(w) \geq 0$ then $\nabla J(w*) = 0$ become a sufficient condition for w* to be the global minimum.

So far, we have studied the functions for the task of regression and classification.
They were all linear functions $y = w^t x$.
And algorithms to learn the functions from data.
(MMSE, Logistic MLE, and SVM dual)

$$\nabla^2 \cdot J(\omega) = \begin{bmatrix} \dfrac{\partial f(w_1, w_2)}{\partial w_1 \partial w_2} & \dfrac{\partial f(w)}{\partial w_1 \partial w} \\ \text{''} & \dfrac{\partial f(w_1, w_2)}{\partial w_2^2} \end{bmatrix}$$

[convexity]

- Logistic Regression MLE:

$$-\frac{\partial \ln \sigma(w^t x)}{\partial w_i} = -\frac{\sigma(w^t x)(1 - \sigma(w^t x))}{\sigma(w^t x)} x_i$$

$$-\frac{\partial^2 \ln \sigma(w^t x)}{\partial w_i \partial w_j} = \sigma(w^t x)(1 - \sigma(w^t x)) x_i x_j \quad > 0$$

At a data point $x$ and current $w$

$$P(t|w) = \prod_{n=1}^{} \sigma(w^t x_n)^{t_n} (1 - \sigma(w^t x_n))^{1 - t_n}$$

$$\arg\min \quad J(\vec{w}) = -\ln P(t|w) = \sum_{n=1}^{N} -t_n \boxed{\ln \sigma(w^t x_n)} - (1 - t_n) \boxed{\ln (1 - \sigma(w^t x_n))}$$

$$\nabla f(x^*) = 0$$

They are convex.

So far, we have studied the functions for the task of regression and classification.
They were all linear functions $y = w^t x$.
And algorithms to learn the functions from data.
(MMSE, Logistic MLE, and SVM dual)

- SVM dual:

$$\begin{bmatrix} t_1 & 0 & \dots & 0 \\ 0 & t_2 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & t_N \end{bmatrix} \cdot \Phi \cdot \Phi^t \cdot \begin{bmatrix} t_1 & 0 & \dots & 0 \\ 0 & t_2 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & t_N \end{bmatrix}$$

$$\lambda *^N_{n=1} = \arg\max_{\lambda*} \sum_{n=1}^{N} \lambda *_n - \frac{1}{2} \sum_{n=1}^{N} \sum_{m=1}^{N} \lambda_n^* \lambda_m^* \cdot t_n \cdot t_m \cdot \kappa(x_n, x_m)$$

Semipositive Definite!

$$\text{subject to} \quad 0 \leq \lambda_n^* \leq C$$

$$\sum_{n=1}^{N} \lambda_n^* \cdot t_n = 0$$

$\nabla f(x*) =$

$-\frac{1}{2}[\lambda \cdots ][\begin{bmatrix} \vdots \end{bmatrix}]$
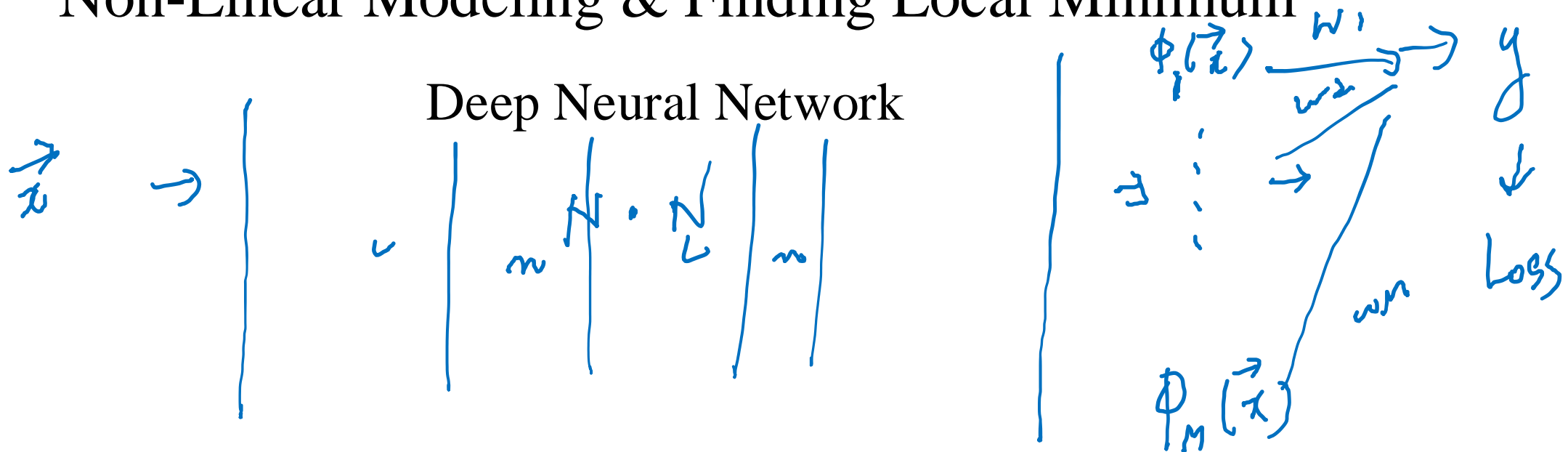
For all linear regression, logistic regression, SVM,

$$\nabla J(w*) = 0$$ is a sufficient condition for the global optimum $w *$.

However, why the methods to find were different?

- Normal Equation,
- Stochastic Gradient Descent,
- SMO (Sequential Minimum Optimization)

$\vec{x}$
data

$\phi_1(\vec{x})$ $w_1$ $\to y$
$\phi_2(\vec{x})$ $w_2$
$\phi_3(\vec{x})$
$\vdots$
$\phi_M(\vec{x})$ $w_M$

Loss

# Non-Linear Modeling & Finding Local Minimum

## Deep Neural Network

$\vec{x}$ $\to$

$\phi_1(\vec{x})$ $w_1$ $\to y$
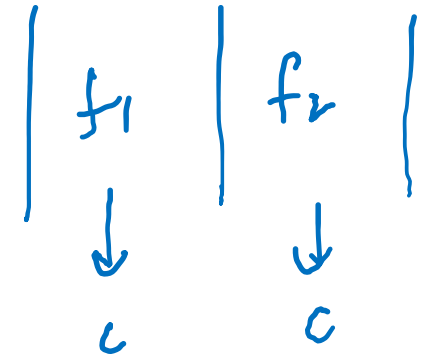$w_2$
$\vdots$
$\phi_M(\vec{x})$ $w_M$ Loss

Deep Neural Networks (DNNs) also learn the heuristic functions $y = f(x)$ for tasks such as classification and regression.
Common loss functions like Mean Squared Error (MSE), Maximum Likelihood Estimation (MLE) for Logistic Regression, and Softmax for multiclass classification are compatible with DNNs. However, the learned function is no longer linear; instead, the network automatically learns a complex feature space through its layers, up to the last hidden layer. The function is a non-linear.

For the non linearity and layered structure
the objective function for DNN is not convex any longer.
Why?

For the non linearity and layered structure,
the objective function for DNN is not convex any longer.
Why?

+ activation functions often result in non-convex functions.
+ A composite function of convex functions is not always convex.

ex) $f(x) = \frac{1}{\sqrt{x}}$ $g(x) = \frac{1}{x}$ $g(f(x)) = ?$ $= \sqrt{x}$

$$\left| \begin{array}{c} f_1 \end{array} \right| \begin{array}{c} f_2 \end{array} \right|$$

$\downarrow$ $\downarrow$

$c$ $c$

$f_2(f_1(x))$

$f$

$1/\sqrt{x}$

When $J(W)$ is non-convex

$\nabla f(x*) = 0$ is a necessary condition for $x*$ to be a local minimum. It has multiple local minimum. Hence, hard to find a global minimum, especially when input to the function is multidimensional.

$$J(w) \downarrow \quad \widehat{w}$$

$$\downarrow$$

$$\nabla J(w^*) = 0$$

Neural Networks are usually trained
by using iterative and gradient based algorithms.

- $w_0$
- $w_1 = w_0 - \eta \nabla J(w_t)$

$$W_{t+1} = W_t - \eta \nabla J(W_t)$$

"Stochastic Method"

How can we compute the gradient $\nabla J(W)$?

$$\nabla_w J(w) = \begin{bmatrix} \dfrac{\partial J(w)}{\partial w_1} \\[2mm] \dfrac{\partial J(w)}{\partial w_2} \\[1mm] ... \\[1mm] \dfrac{\partial J(w)}{\partial w_{M-1}} \\[2mm] \dfrac{\partial J(w)}{\partial w_M} \end{bmatrix}$$

# of parameters to learn

- AlexNet : 62.3M
- VGG-16: 138 M
- Chat-GPT2: 1.5 B

# Computing Gradient $\nabla J(W)$
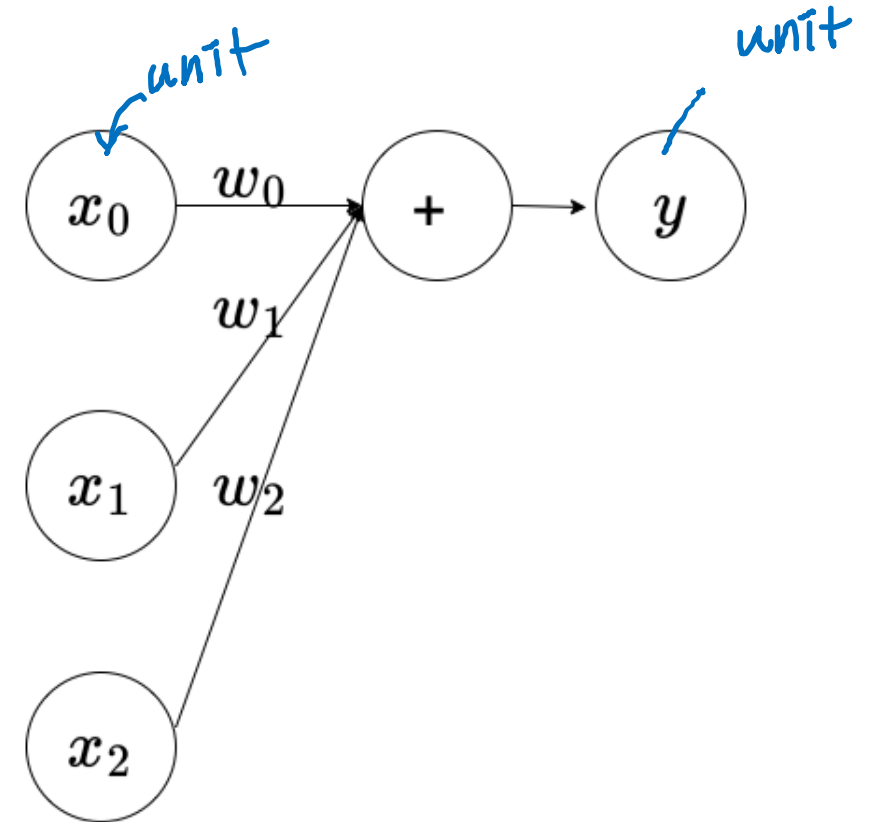## Backpropagation Algorithm

A function $f \colon X \to Y$
A function $X$ from Y where $X$ is the domain and $Y$ is codomain.
$f(x) = y$, $y$ is the image of $x$.

Ex) One Elementary Function of DNN
$$y = w^t x$$

A function is called multivariable
if its input is made up of multiple members.

A function $f\colon\ X \subseteq \mathcal{R}^M \to Y$
A function $X$ from $Y$, we say $X$ is the domain and $Y$ is codomain.
$f(x)\ =\ y$, $y$ is the image of $x$
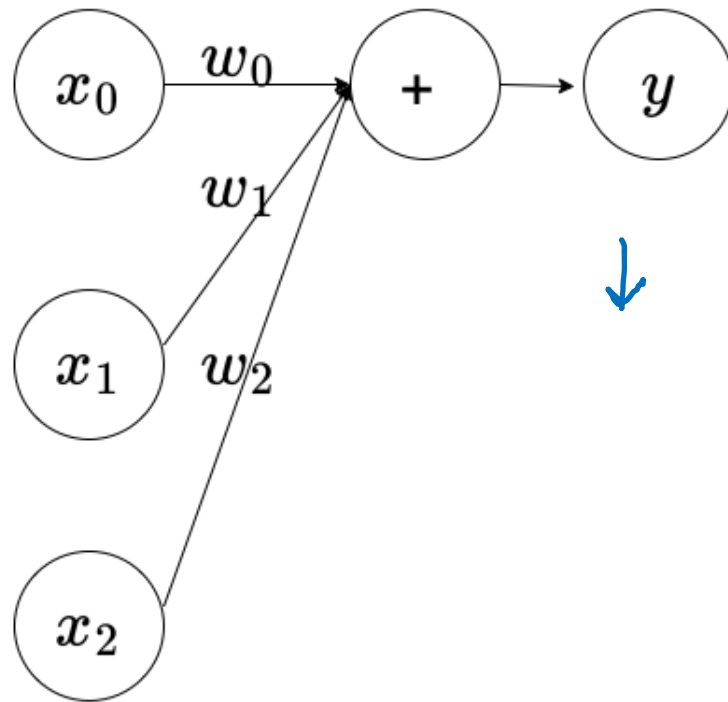The range (or image) of $f$ is the set of all images of $A$.

If the output function consists of multiple values,
it is also be called <u>multivariable</u>,
but these one are called <u>vector valued functions</u>.

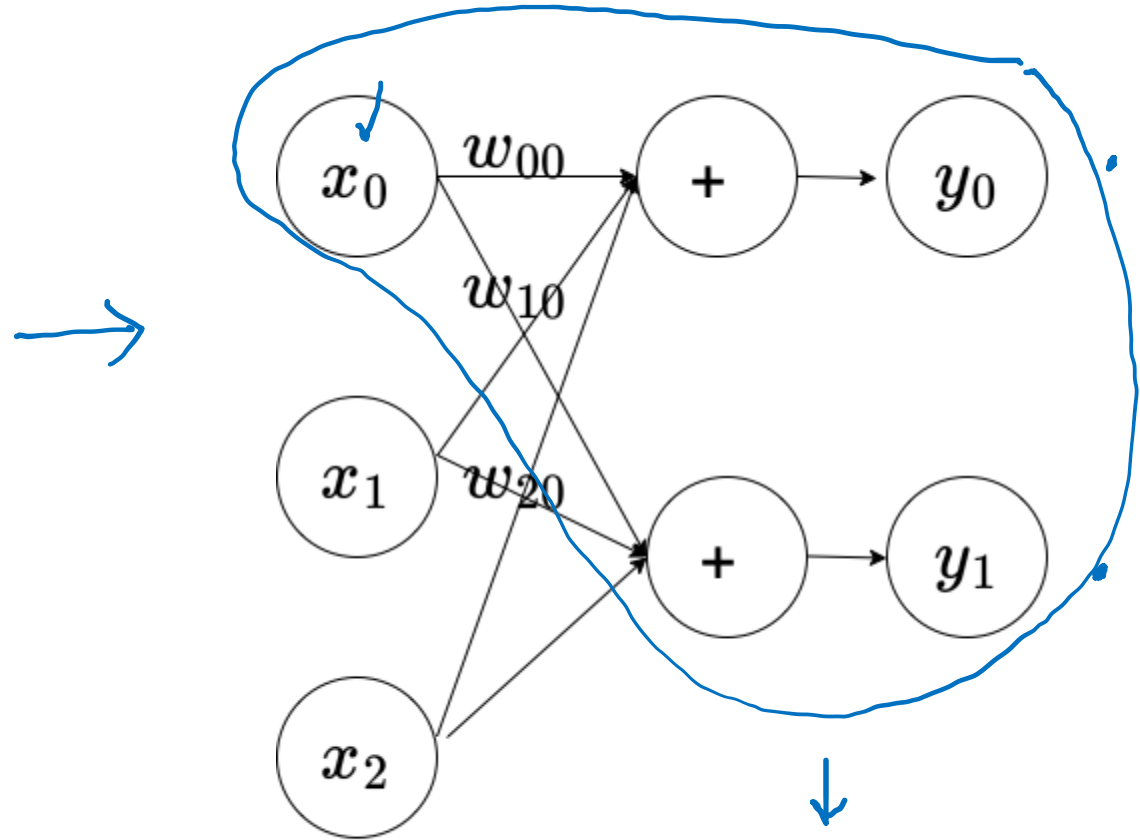A function $f:\ X \subseteq \mathcal{R} \to Y \subseteq \mathcal{R}^N$
A function $X$ from $Y$, we say $X$ is the domain and $Y$ is codomain.
$f(x)\ =\ y$, $y$ is the image of $x$
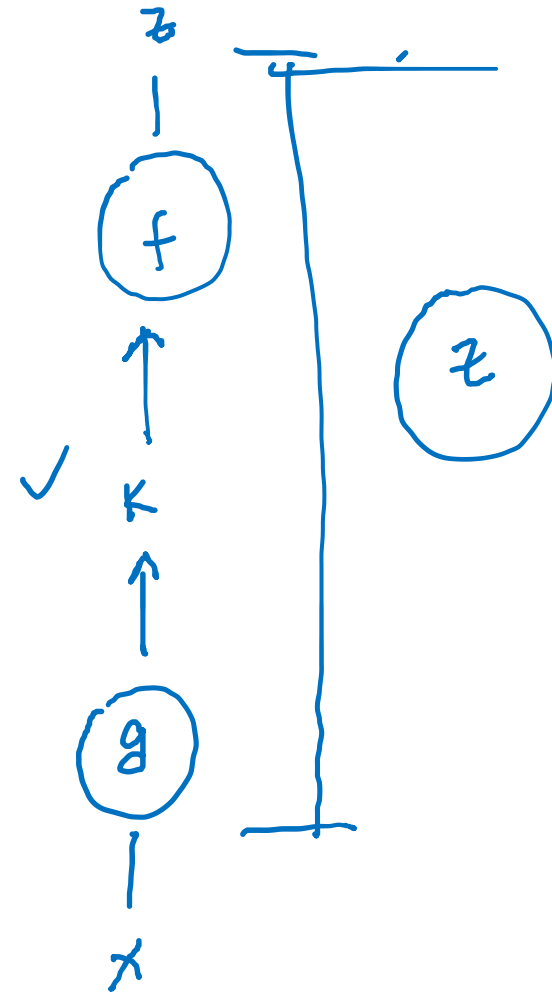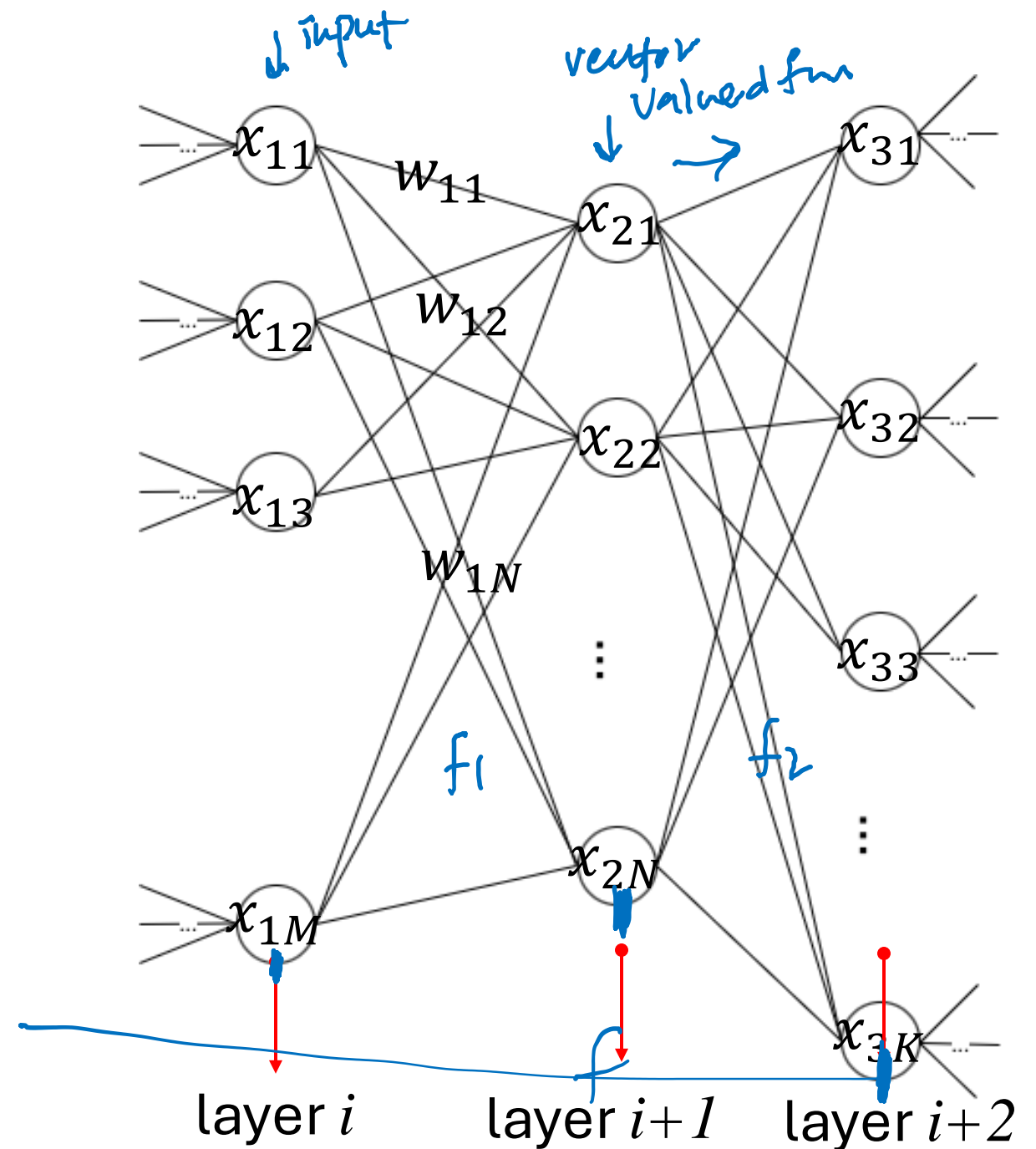The range  (or image) of f is the set of all images of $A$.

Multivariable Function

Multivariable and Vector Valued Function

A composite function $f \circ g \colon X \to Y$ :
where $g \colon X \to Z$ and $f \colon Z \to Y$

$$f = f_2\left(\left(f_1(x)\right)\right)$$

- DNN is <u>a large composite function</u> consisting of many elementary Vector Valued Functions.

- Each vector-valued function becomes one layer.

- The chain structure gives the depth.



$d$ input

vector valued fn

$x_{11}$

$w_{11}$

$x_{21}$

$x_{31}$

$x_{12}$

$w_{12}$

$x_{22}$

$x_{32}$

$x_{13}$

$w_{1N}$

$x_{33}$

$f_1$

$f_2$

$x_{2N}$

$x_{1M}$

$x_{3K}$

layer $i$    layer $i{+}1$    layer $i{+}2$

# Derivative & Partial Derivative

The derivative at point $x_0$ of a scalar function $f$ :

$$f'(x_0) = \frac{df(x_0)}{dx} = \lim_{h \to 0} \frac{\overbrace{f(x_0 + h) - f(x_0)}^{\Delta y}}{\underbrace{h}_{\Delta x}}$$

The partial derivative of a function with two variables $f(x, y)$ r.t $x$ at point $(x_0, y_0)$.

$$f_x(x_0, y_0) = \frac{\partial f}{\partial x}(x_0, y_0) = \lim_{h \to 0} \frac{f(x_0 + h, y_0) - f(x_0, y_0)}{h}$$

$f'(x)$

$$\nabla_w J(w) = \begin{bmatrix} \dfrac{\partial J(w)}{\partial w_1} \\[2mm] \dfrac{\partial J(w)}{\partial w_2} \\[1mm] \dots \\[1mm] \dfrac{\partial J(w)}{\partial w_{M-1}} \\[2mm] \dfrac{\partial J(w)}{\partial w_M} \end{bmatrix}$$

DNN's $J(w)$ represents a multivariate function, so we need to compute the partial derivatives.

The second or higher order derivatives can be defined.
Some gradient descent algorithm use the second order derivative information,
but today we are going to focus on computing <span style="color:red">first order derivatives</span>.

$$f''(x_0) = \frac{\mathrm{d}^2 f(x_0)}{\mathrm{d}x^2}$$

Scalar function

$$f_{xx}(x_0, y_0) = \frac{\partial^2 f(x_0, y_0)}{\partial x^2}$$

The function multiple variable

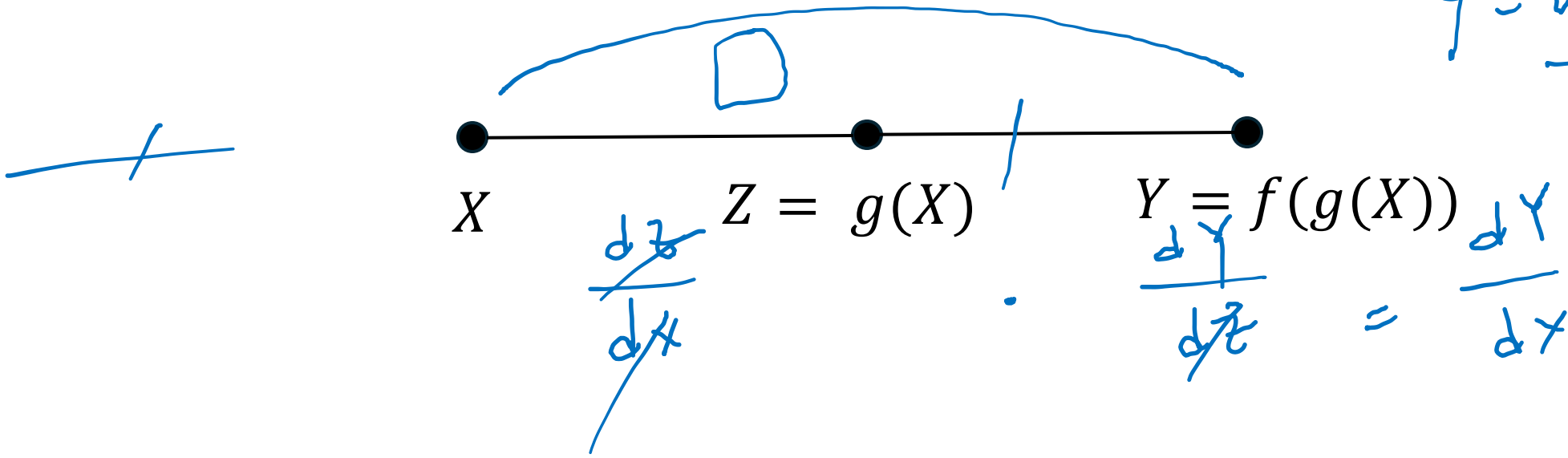$$f_{xy}(x_0, y_0) = \frac{\partial^2 f(x_0, y_0)}{\partial x \partial y}$$

$$g(x) = x$$

$$f(z) = z^2$$

$$f(g(x)) = \boxed{x^2}$$

# Chain Rule

## Computing the Derivative of Composite Functions

$$Y = w(x)$$



$$X \qquad Z = g(X) \qquad Y = f(g(X))$$
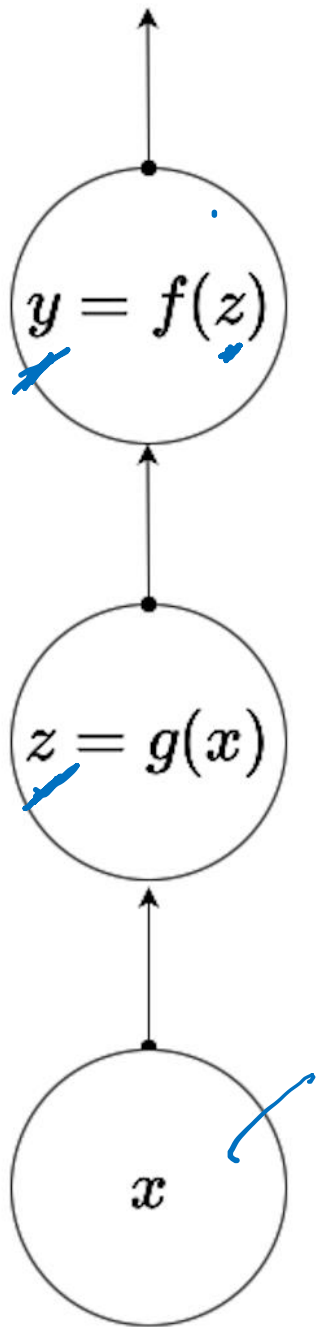
$$\frac{dz}{dx} \qquad \frac{dY}{dz} = \frac{dY}{dx}$$

# Chain Rule for Scalar Functions
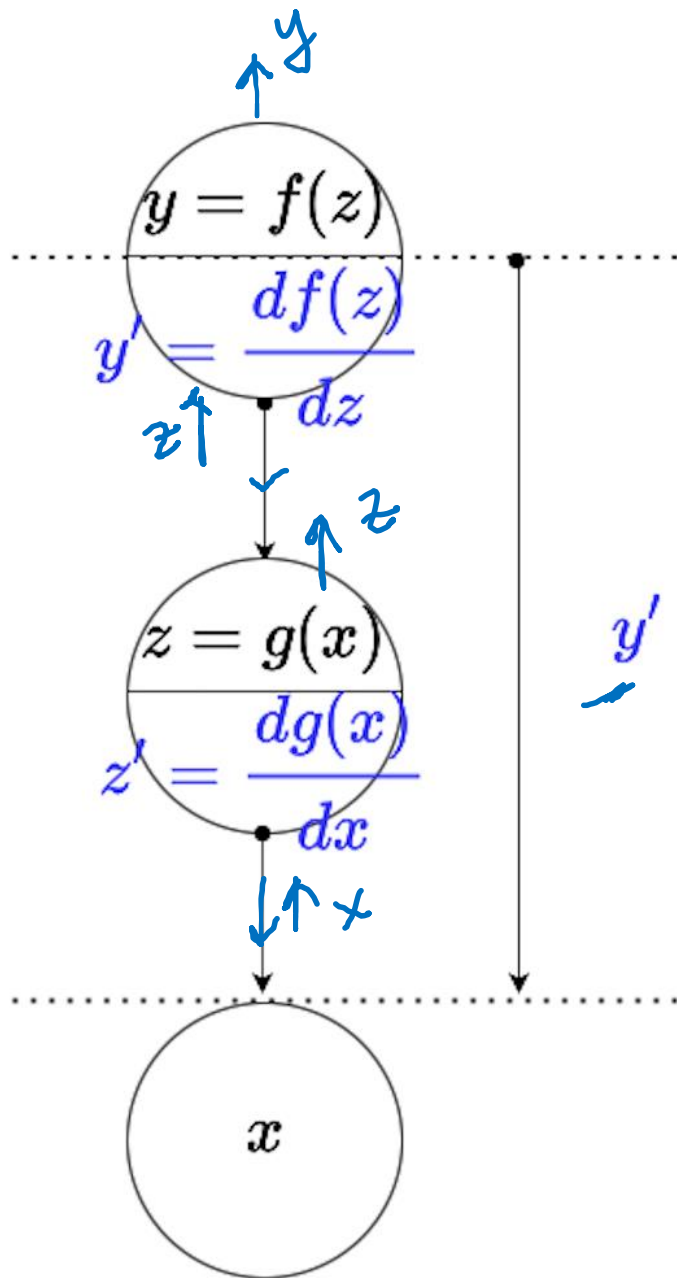
Let $z = g(x)$ and $y = f(z)$,

Then $\dfrac{d(f \circ g\,(x))}{dx} = \dfrac{dy}{dx} = \dfrac{dz}{dx} \cdot \dfrac{dy}{dz}$

$y = f(z)$

$z = g(x)$

$x$

unit.

- The feedforward direction to compute a composite function $f \circ g$

28

- The backward direction to compute the derivate of composite function $\dfrac{df \circ g\,(x)}{dx}$

$$y' = \frac{df(g(x))}{dx} = \frac{df(z)}{dz} \cdot \frac{dg(x)}{dx}$$

input

ouput

$$\frac{d\,Loss}{dw} = \frac{d\,Loss}{dx} \cdot \left(\frac{dx}{dw}\right)$$

- Computing Gradients in a Neural Net:

Q: How can we compute the $\frac{dLoss}{dw}$ by using $\frac{dLoss}{dx}$?

**Diagram (left side):**

$Loss = y$

$\frac{\partial y}{\partial y} = 1$

$y = f(z)$

$y' = \frac{df(z)}{dz}$

$\frac{df(z)}{dz} \cdot 1$

$z = g(x)$

$z' = \frac{dg(x)}{dx}$

$\frac{df(z)}{dz} \cdot 1$

$\frac{dg(x)}{dx} \cdot$

$x$

$w$

$d$
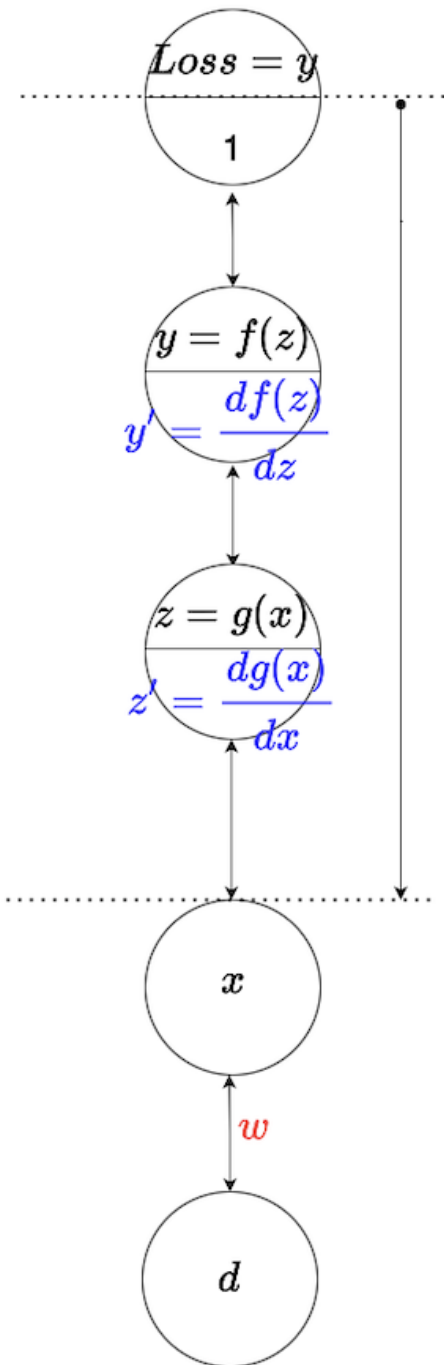
$$\frac{dLoss(x)}{dx} = 1 \cdot \frac{df(z)}{dz} \cdot \frac{dg(x)}{dx}$$

$$\frac{d\,Loss}{dx} = \frac{d\,Loss}{dw} \cdot \frac{dw}{dx}$$

$x = w \cdot d$

$\frac{dx}{dw} = d$

- Computing Gradients in a Neural Net:

$$\frac{dLoss}{dw} = \frac{dLoss}{dx} \cdot \frac{dx}{dw} = \frac{dLoss}{dx} \cdot d$$

$$\frac{dLoss(x)}{dx} = 1 \cdot \frac{df(z)}{dz} \cdot \frac{dg(x)}{dx}$$

# Chain Rule for Multivariable Functions

$$\frac{dw}{dt} = \frac{\partial x}{dt} \cdot \frac{\partial f}{\partial x} \boxed{t} \quad \frac{dy}{dt} \cdot \frac{\partial f}{\partial y}$$

Let $w = f(x, y),$ where $x = x(t)$ and $y = y(t)$

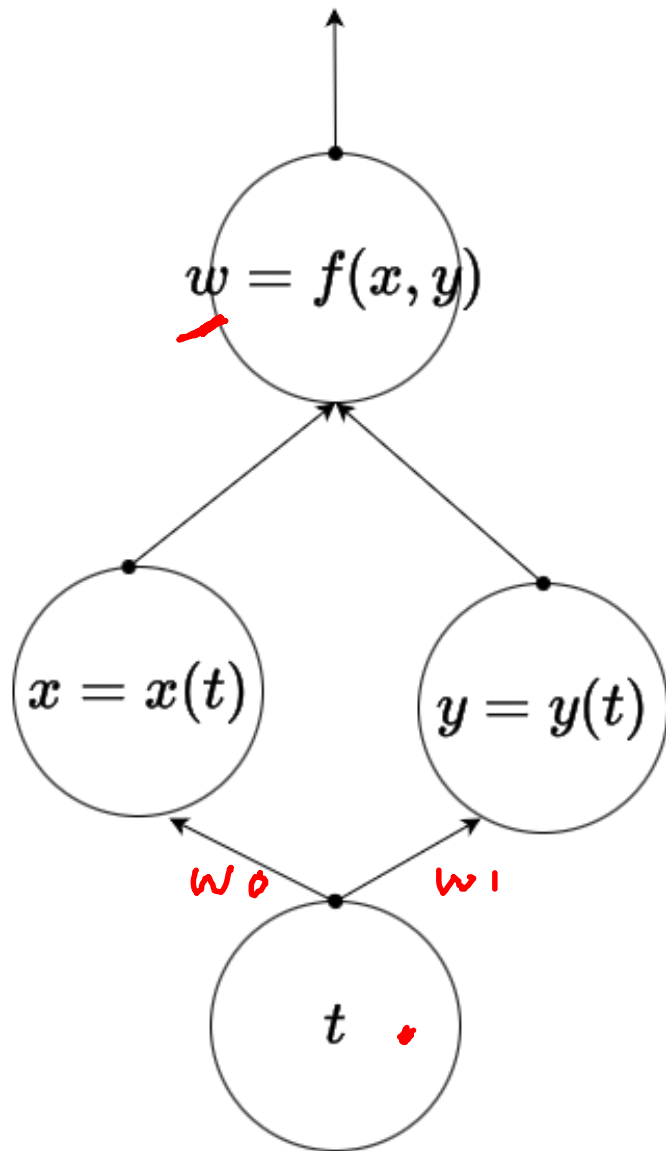Then $\dfrac{df(x,y)}{dt} = \dfrac{dw}{dt} = \dfrac{\partial w}{\partial x} \cdot \dfrac{dx}{dt} + \dfrac{\partial w}{\partial y} \cdot \dfrac{dy}{dt}$

$\dfrac{dx}{dt}$    $\dfrac{\partial f}{\partial x}$

$\dfrac{dy}{dt}$    $\dfrac{\partial f}{\partial y}$

$t$        $x = x(t)$        $w = f(x, y)$

$y = y(t)$

- The feedforward direction to compute a composite function
$$w(t) = f(x(t), y(t))$$

- The backward direction to compute the derivate of composite function $\dfrac{dw}{dt}$

$$w = f(x,y)$$

$$\frac{\partial f(x,y)}{\partial x} \qquad \frac{\partial f(x,y)}{\partial y}$$

$$x = x(t) \qquad y = y(t)$$

$$\frac{dx(t)}{dt} \qquad \frac{dy(t)}{dt}$$

$$t$$

$$\frac{\partial f(x,y)}{\partial x} \cdot \frac{dx(t)}{dt} \qquad\qquad \frac{\partial f(x,y)}{\partial y} \cdot \frac{dy(t)}{dt}$$

$$\frac{dw}{dt} = \frac{\partial f(x,y)}{\partial x} \cdot \frac{dx(t)}{dt} + \frac{\partial f(x,y)}{\partial y} \cdot \frac{dy(t)}{dt}$$

$$\frac{Loss = y}{1}$$

$$w = f(x, y)$$

$$\frac{\partial f(x,y)}{\partial x} \qquad \frac{\partial f(x,y)}{\partial y}$$

$$\frac{\partial f(x,y)}{\partial x} \cdot \frac{dx(t)}{dt} \qquad \frac{\partial f(x,y)}{\partial y} \cdot \frac{dy(t)}{dt}$$

$$\frac{x = x(t)}{\frac{dx(t)}{dt}} \qquad \frac{y = y(t)}{\frac{dy(t)}{dt}}$$

$t$

$w$

$d$

$$\frac{dw}{dt} = 1 \cdot \underbrace{\frac{\partial f(x,y)}{\partial x} \cdot \frac{dx(t)}{dt}}_{\text{①}} + \underbrace{\frac{\partial f(x,y)}{\partial y} \cdot \frac{dy(t)}{dt}}_{\text{②}}$$

Loss

$\frac{dL}{dt}$

- Computing Gradients in a Neural Net:

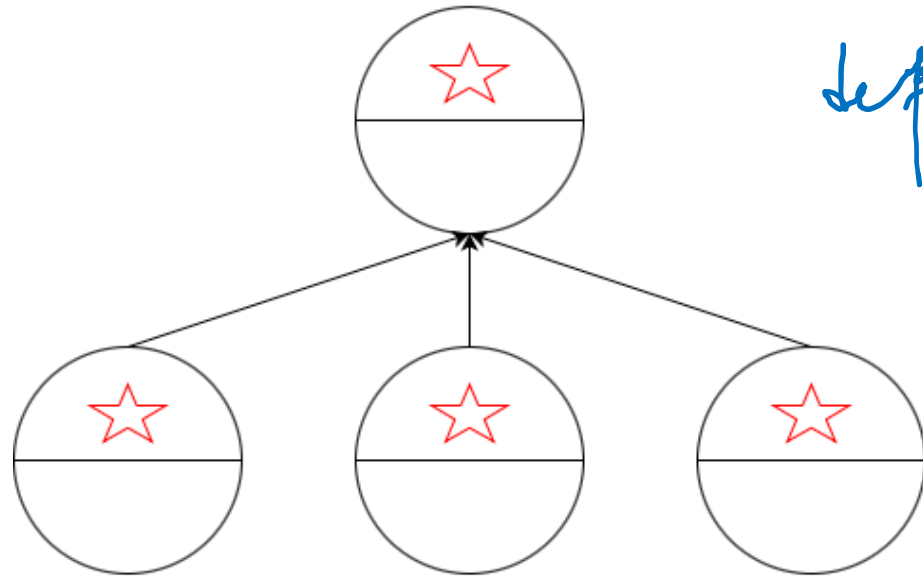  Q: How can we compute the $\frac{dLoss}{dw}$ by using $\frac{dLoss}{dt}$?

35

At a Certain $W$, $[W_{t+1} = W_t - \eta \nabla J(W_t)]$

The Derivative of Loss at Each Node is Updated through Backpropagation.



$$\frac{dLoss}{dx_{00}} = \frac{dx_{10}}{dx_{00}}\frac{\partial Loss}{\partial x_{10}} + \frac{dx_{11}}{dx_{00}}\frac{\partial Loss}{\partial x_{11}} + \frac{dx_{12}}{dx_{00}}\frac{\partial Loss}{\partial x_{12}}$$

class unit ( ) { memory,

def feedforward ( )
      operation ( input , output )

def feedback ( )

- location
  cal information

- input , local inf

- out

Feedforward Update
(For Inference)
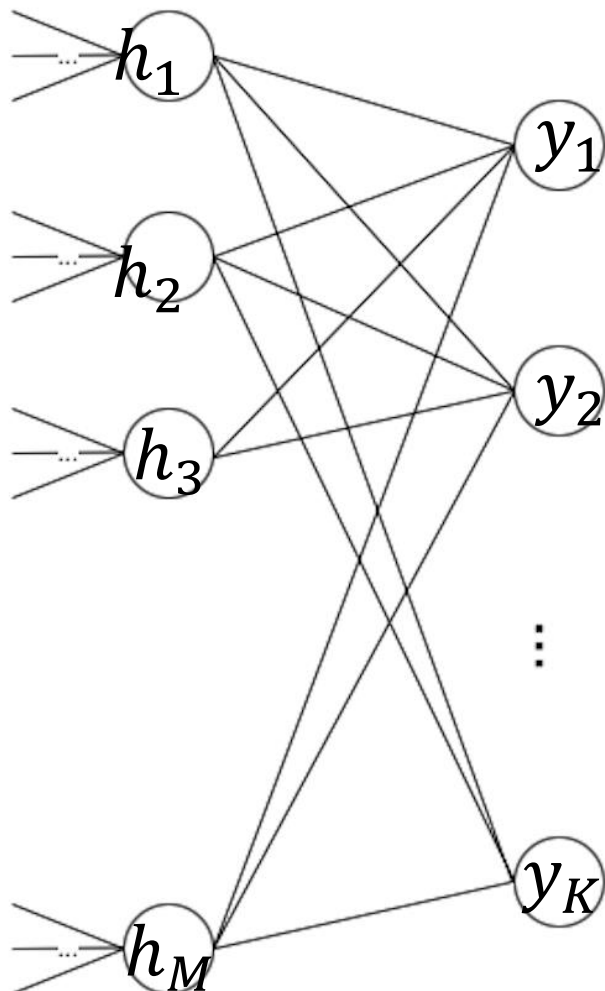
Feedbackward Update
(For Computing Gradient)

# Output Layers of DNN and Gradient Descent

# [1] MMSE Loss for Regression

# [1] MMSE Loss for Regression



In regression modeling,

(1) $\vec{y} = f(\vec{x}) + \varepsilon$ where $\varepsilon \sim \mathcal{N}(0, \Sigma)$

$$then \ y \sim \mathcal{N}(f(x), \Sigma)$$
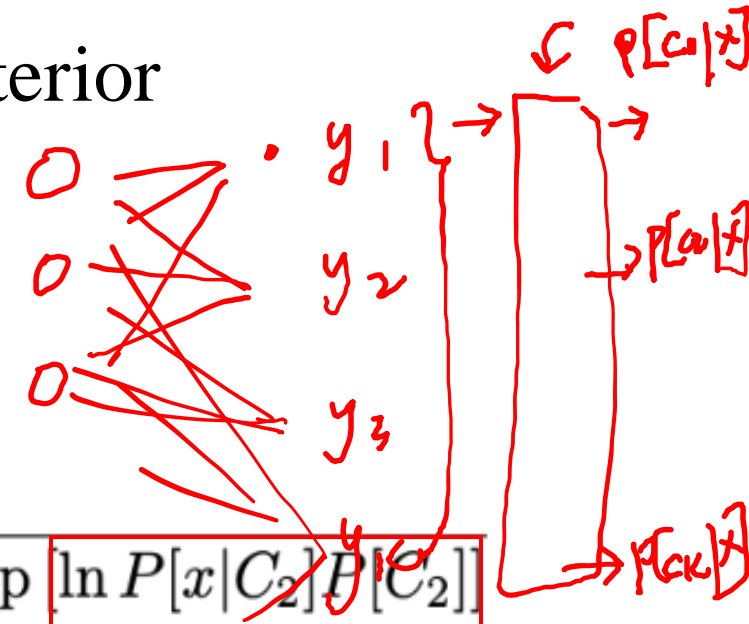
(2) For MLE,

Take a NLL (negative log likelihood)

# [2] Softmax for Multiple Classification

K−class → #K discriminant

# SoftMax Representation of Multiclass Posterior

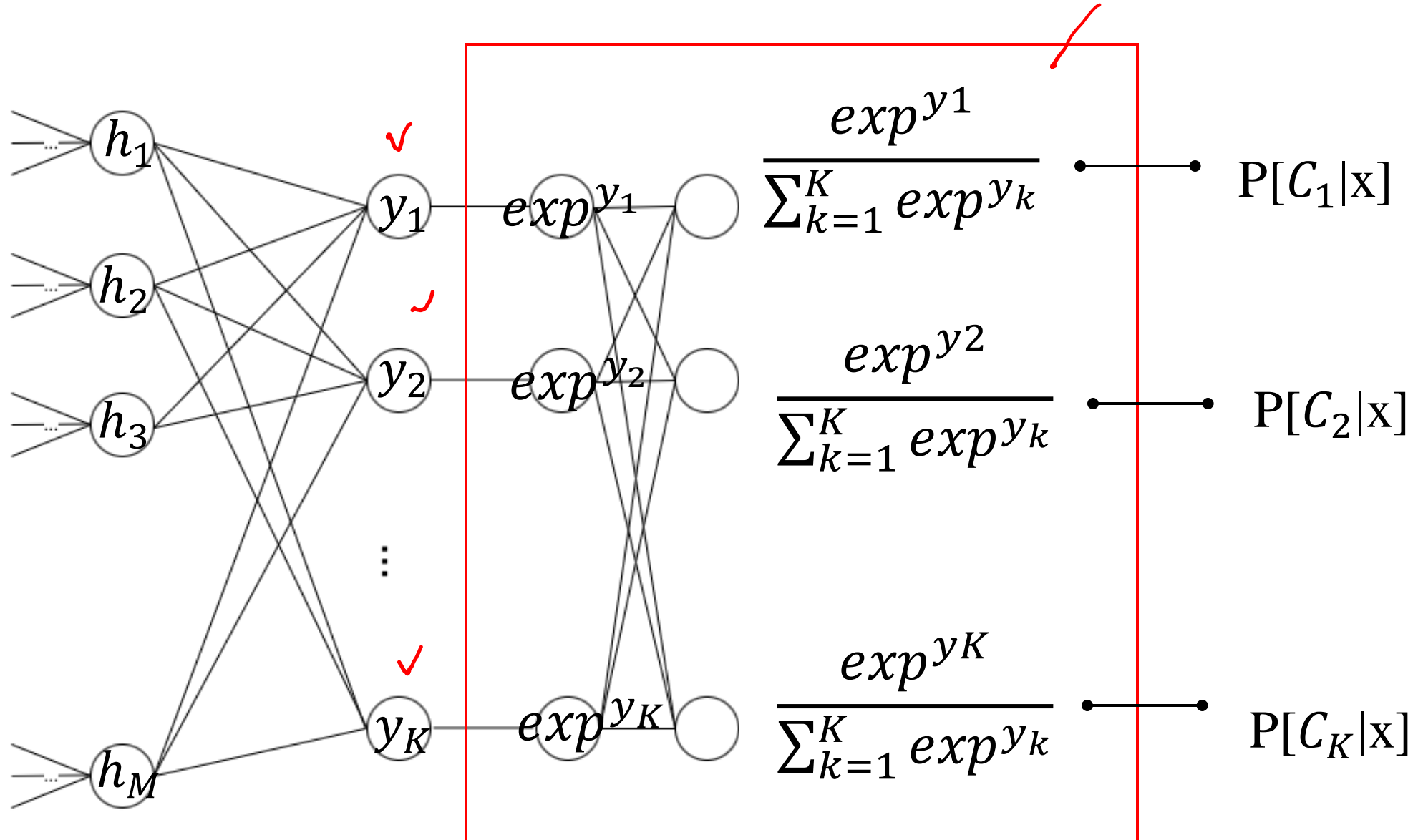$$P[C_0|x] = \frac{P[x|C_0]P[C_0]}{P[x|C_0]P[C_0] + P[x|C_1]P[C_1] + P[x|C_2]P[C_2]}$$

$$P[C_0|x] = \frac{\exp[\ln P[x|C_0]P[C_0]]}{\exp[\ln P[x|C_0]P[C_0]] + \exp[\ln P[x|C_1]P[C_1]] + \exp[\ln P[x|C_2]P[C_2]]}$$
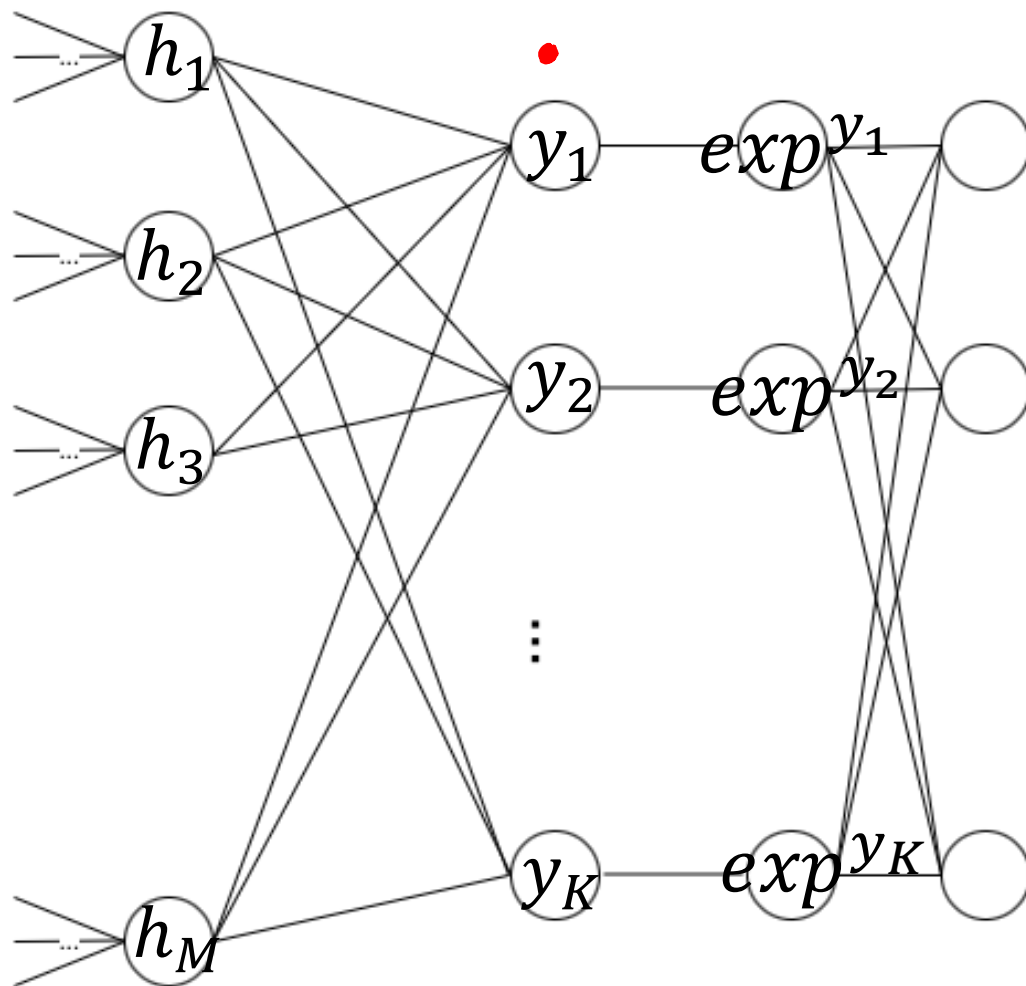
$P[c_1|x]$

$P[c_2|x]$

$P[c_K|x]$

$y_1$
$y_2$
$y_3$

+ K= 3 discriminant functions

The posterior will be designed with softmax & $K = 3$ linear function of $x$

& $K = 3$ linear function of $\phi(x)$

# [2] SoftMax Ouput for Multiclass Classification



$$\frac{exp^{y1}}{\sum_{k=1}^{K} exp^{y_k}} \quad \text{P}[C_1|\text{x}]$$

$$\frac{exp^{y2}}{\sum_{k=1}^{K} exp^{y_k}} \quad \text{P}[C_2|\text{x}]$$

$$\frac{exp^{yK}}{\sum_{k=1}^{K} exp^{y_k}} \quad \text{P}[C_K|\text{x}]$$

# [3] SoftMax Unit and MLE Loss

$$P[T|W,x] = \prod_{n=1}^{N} \prod_{k=1}^{K} P[C_k|W,x_n]^{t_{nk}}$$

$$J(W) = -\ln P[T|W,x] = -\sum_{n=1}^{N} \sum_{k=1}^{K} t_{nk} \ln P[C_k|W,x_n]$$

Once we know $\frac{dJ}{dy_j}$,

we can compute the gradients
of other units by backpropagation.

# Sofmax Derivative r.t $y_j$

$$P[T|W,x] = \prod_{n=1}^{N} \prod_{k=1}^{K} P[C_k|W,x_n]^{t_{nk}}$$

$$J(W) = -\ln P[T|W,x] = -\sum_{n=1}^{N} \sum_{k=1}^{K} t_{nk} \ln P[C_k|W,x_n]$$

$$\frac{\partial J}{\partial y_j} = \sum_{n=1}^{N} \{t_{n1}P[C_j|x_n,W] + ... + t_{nj}(P[C_j|x_n,W]-1) + t_{nK}P[C_j|x_n,W]\}$$

$$\frac{\partial J}{\partial y_j} = \sum_{n=1}^{N} \{P[C_j|x_n,W] - t_{nj}\}$$
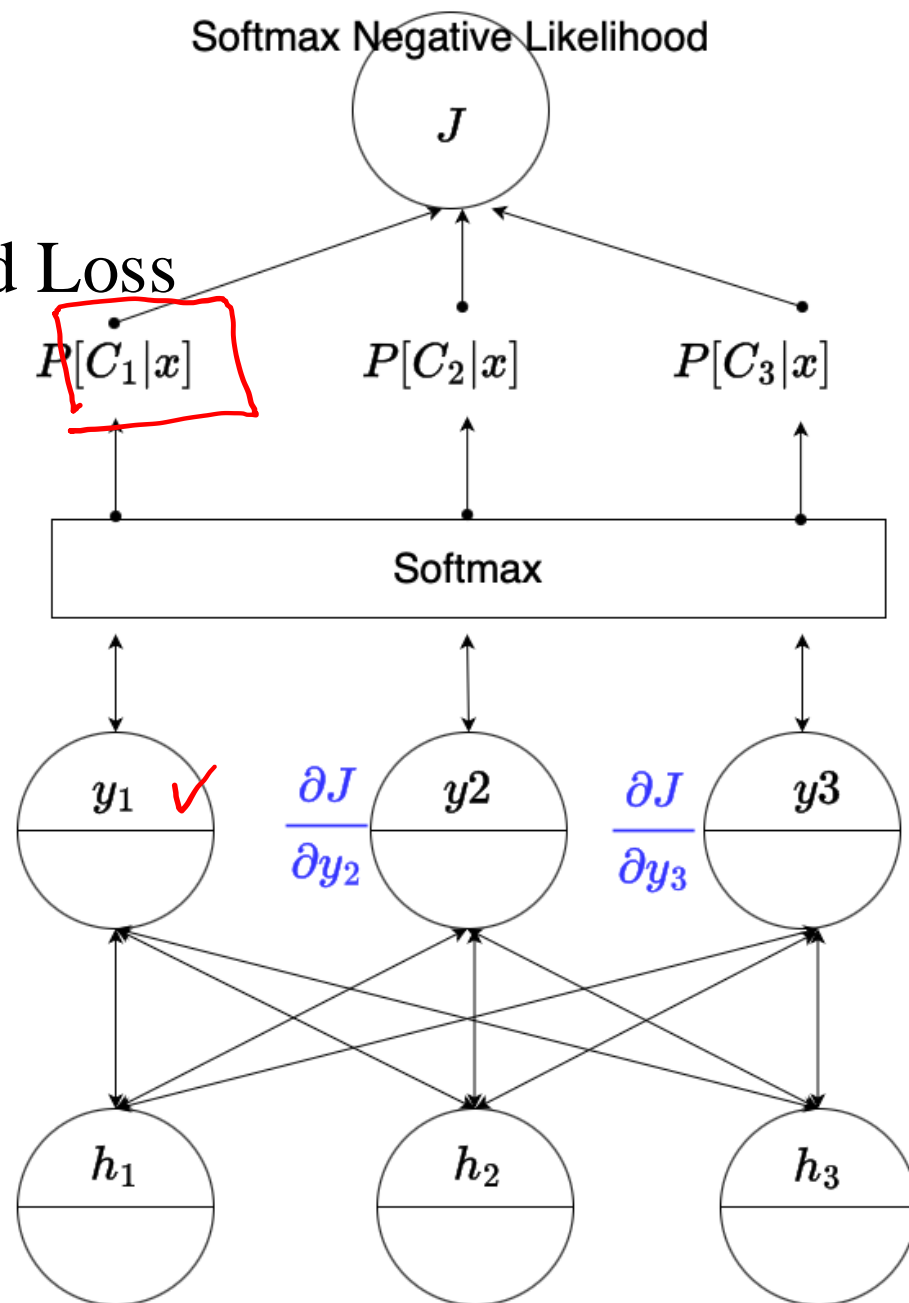
# SoftMax Derivative [Slide Oct. 7th]

$$\checkmark$$

$$\nabla_{W_j} P[C_j|x] = P[C_j|x] \cdot (1 - P[C_j|x])\phi(x)$$

$$\nabla_{W_i} P[C_j|x] = P[C_j|x] \cdot (-P[C_i|x])\phi(x)$$

+ we can compute the P[Cj|x] at current W!

Softmax Negative Likelihood

$$J$$

The derivatives of Softmax Negative Likelihood Loss
: they backpropagated throughout the network
to compute the gradients of the network units.
The parameters connected to the units
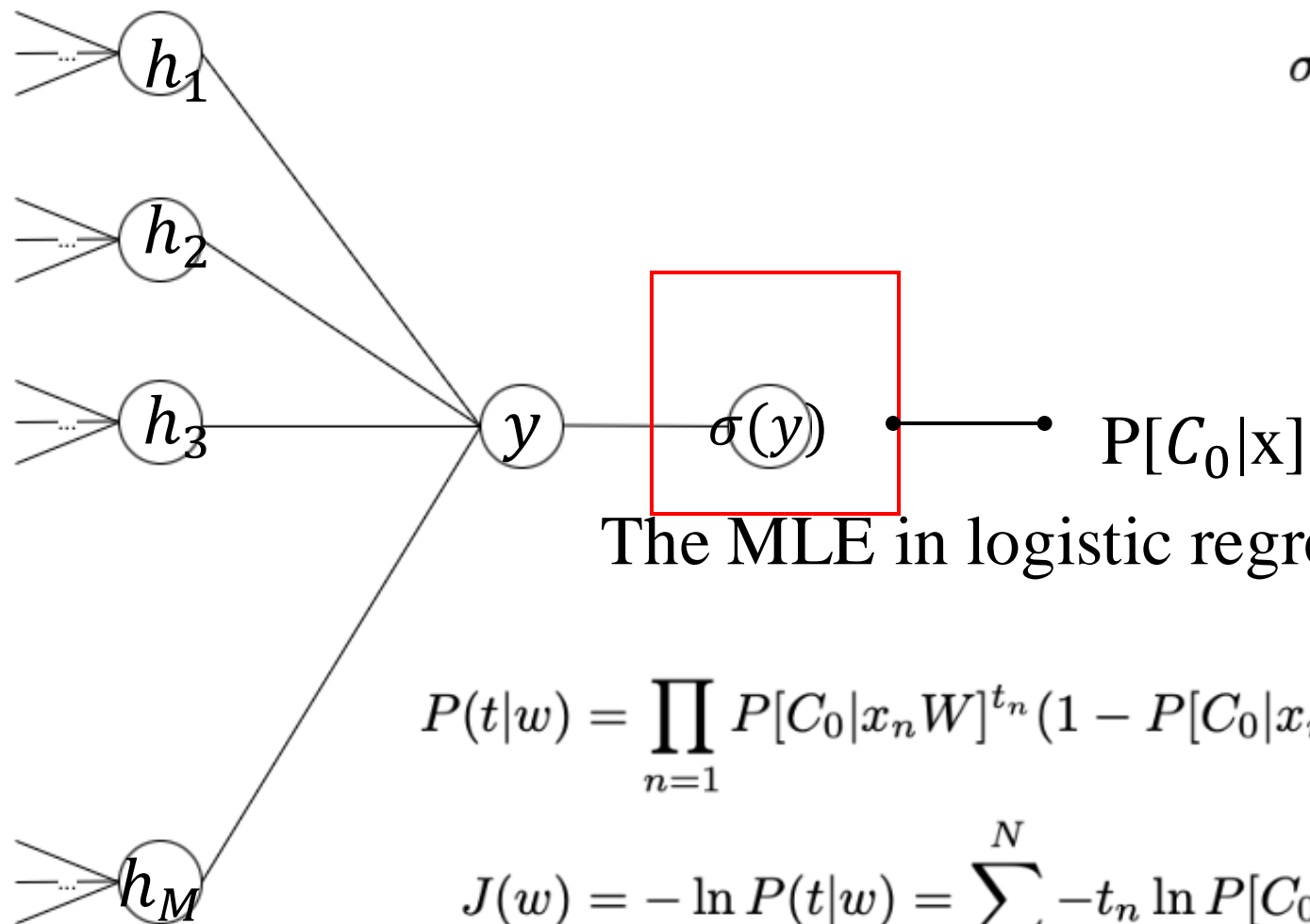are updated based on the updated gradients.

$P[C_1|x]$  $P[C_2|x]$  $P[C_3|x]$

Softmax

$$\frac{\partial J}{\partial y_1} = P[C_1|x] - t1$$

$y_1$   $\frac{\partial J}{\partial y_2}$ $y2$   $\frac{\partial J}{\partial y_3}$ $y3$

$h_1$   $h_2$   $h_3$

# [2] Logistic Regression for Binary Classification

# [3] Sigmoid Unit for Binary Classification



$$\sigma(y) = \frac{1}{1 + \exp{-y}} = \frac{1}{1 + \exp{-w^t h}}$$

$P[C_0|x]$

The MLE in logistic regression modeling,

$$P(t|w) = \prod_{n=1} P[C_0|x_n W]^{t_n} (1 - P[C_0|x_n W])^{1-t_n}$$

KL divergence!
Why?

$$J(w) = -\ln P(t|w) = \sum_{n=1}^{N} -t_n \ln P[C_0|x_n W] - (1 - t_n) \ln (1 - P[C_0|x_n W])$$

# Logistic Regression Derivative r.t $y$

$$J(w) = -\ln P(t|w) = \sum_{n=1}^{N} -t_n \ln P[C_0|x_n W] - (1 - t_n) \ln (1 - P[C_0|x_n W])$$

$$\frac{\partial J}{\partial y} = \sum_{n=1}^{N} -t_n (1 - \sigma(y|x_n)) - (1 - t_n)(-\sigma(y|x_n))$$

$$= \sum_{n=1}^{N} \sigma(y|x_n) - t_n$$

Logistic Negative Likelihood

$J$

The derivatives of Sigmoid Negative Likelihood Loss
: it is backpropagated throughout the network
to compute the gradients of the network units.
The parameters connected to the units
are updated based on the updated gradients.

$P[C_0|x]$

Logistic Sigmoid

$$\frac{\partial J}{\partial y} = P[C_0|x] - t$$

$y$

$h_1$

$h_2$

$h_3$