# Introduction to Deep Learning
## Calculus behind Backpropagation

Lecturer: Dongyu Yao

$x_1$  $w_{1,1}^{(1)}$  $v_1^{(1)}$  $w_{1,1}^{(2)}$  $v_1^{(2)}$  $w_{1,1}^{(3)}$

$x_2$  $v_2^{(1)}$  $v_2^{(2)}$  $f_1$

$\vdots$

$x_d$  $v_{n_1}^{(1)}$  $v_{n_2}^{(2)}$  $f_k$

$w_{n_1,d}^{(1)}$  $w_{n_2,n_1}^{(2)}$  $w_{k,n_2}^{(3)}$

Input
layer

Hidden
layer 1

Hidden
layer 2

Output
layer

# Forward propagation (short notation, with bias nodes)

- For input layer: $v^{(0)} = [x; \ 1]$

- For each hidden layer $l = 1{:}L - 1$

$$z^{(l)} = W^{(l)} v^{(l-1)} \quad \text{and} \quad v^{(l)} = [\varphi(z^{(l)}); \ 1]$$

- For output layer:

$$f = W^{(L)} v^{(L-1)}$$

- Predict $y = f$ for regression, or $y = \arg \max_j f_j$ for classification

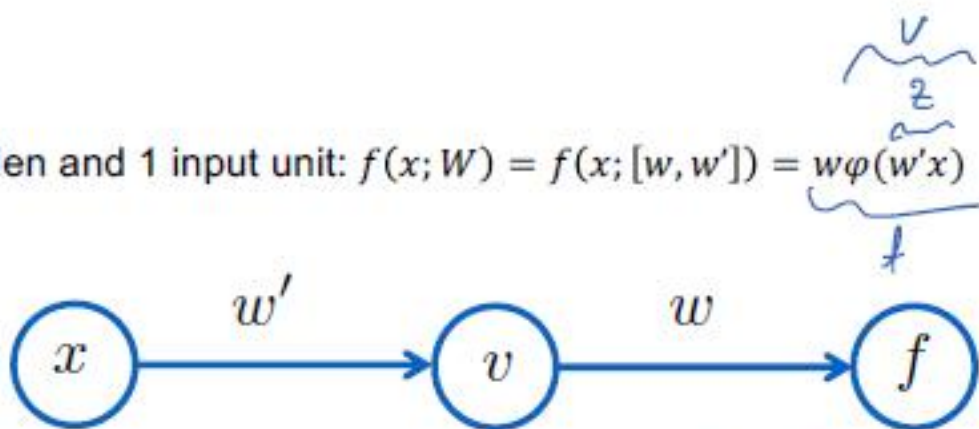# Computing the gradient

In order to apply SGD, need to compute

$$\nabla_{W}\ell(W; x, y)$$

I.e., for each weight between any two connected units $i$ and $j$ on layer $l$ need

$$\frac{\partial}{\partial w_{i,j}^{(l)}}\ell(W; x, y)$$

# Simple example

ANN with 1 output, 1 hidden and 1 input unit: $f(x; W) = f(x; [w, w']) = w\varphi(\overbrace{w'x})$

$\overset{v}{\overset{z}{\underset{f}{}}}$



Loss: $\ell(W; x, y)$, e.g. $\ell(V; x, y) = [y - f(x; W)]^2 =: \ell_y(f)$

$\dfrac{\partial \ell}{\partial w} = \dfrac{\partial \ell}{\partial f} \dfrac{\partial f}{\partial w} = \ell'_y(f) \cdot v$ ← Computed during forward pass

$\underset{\delta^1}{}$ $\underset{\delta}{}$

$\dfrac{\partial \ell}{\partial w'} = \underset{\delta}{\dfrac{\partial \ell}{\partial f}} \cdot \underset{w}{\dfrac{\partial f}{\partial v}} \cdot \underset{\varphi'(z)}{\dfrac{\partial v}{\partial z}} \underset{x}{\dfrac{\partial z}{\partial w'}}$

Want:

$\nabla_W \ell(V; x, y)$

$= \left[ \dfrac{\partial \ell}{\partial w}, \dfrac{\partial \ell}{\partial w'} \right]^T$

$z = w' \cdot x$

$v = \varphi(z)$

$f = w \cdot v$

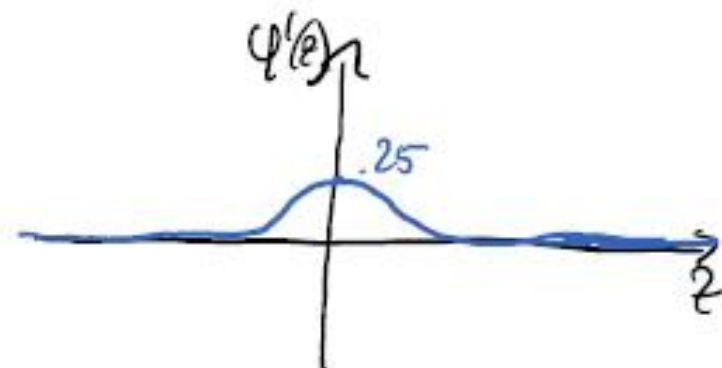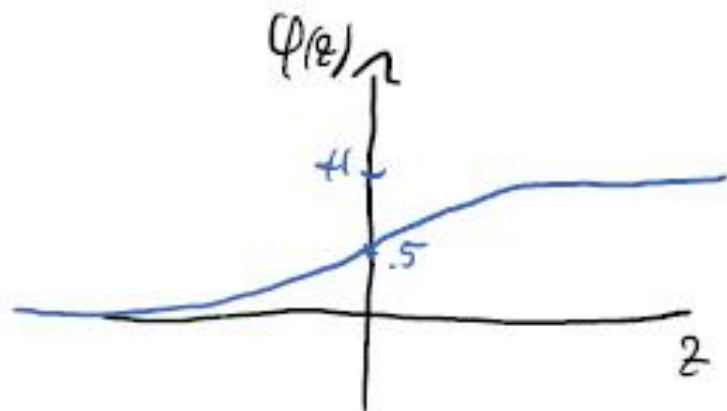$\ell'_y(f) = 2(f - y) =: \delta$

# Derivatives of activation functions

$$\varphi(z)\nearrow$$

Sigmoid:
$$\varphi(z) = \frac{1}{1+\exp(-z)}$$

$$\varphi'(z) = \frac{-1}{(1+e^{-z})^2} \cdot e^{-z} \cdot (-1) = \frac{e^{-z}}{(1+e^{-z})^2}$$

$$= \frac{e^{-z}}{(1+e^{-z})} \cdot \frac{1}{(1+e^{-z})} = \underbrace{(1-\varphi(z))}_{(1-v)} \underbrace{\varphi(z)}_{v}$$

+ $\varphi'(z)$ is easy to compute given $\varphi(z) = v$

+ $\varphi'(z) \neq 0 \; \forall z$, defined everywhere

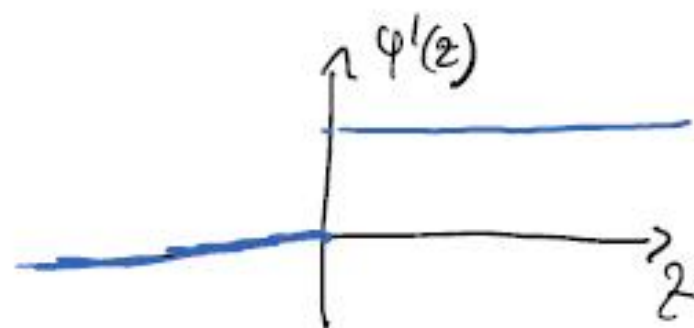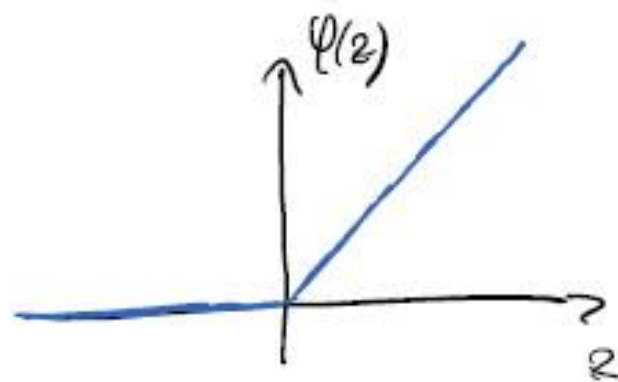− $\varphi'(z) \approx 0 \; \forall z$ (except $z \approx 0$)

# Derivatives of activation functions

Rectified Linear Unit (ReLU):     $\varphi(z) = \max(z, 0)$

$$\varphi'(z) = \begin{cases} 1 & \text{if } z > 0 \\ 0 & \text{if } z < 0 \end{cases}$$

- not diff'able at $z = 0$
  doesn't matter much in practice;
  e.g. simply define $\varphi'(0) := 0$

+ even easier to compute than sigmoid

+ $\not= 0$    $\forall z > 0$

# Recall: Jacobians & multivariate chain rule

For a function $f : \mathbb{R}^n \to \mathbb{R}^m$ its Jacobian at x is given by

$$\frac{d}{d\boldsymbol{x}} f(x) = \begin{pmatrix} \dfrac{\partial f_1}{\partial x_1} & \cdots & \dfrac{\partial f_1}{\partial x_n} \\ \vdots & & \vdots \\ \dfrac{\partial f_m}{\partial x_1} & \cdots & \dfrac{\partial f_m}{\partial x_n} \end{pmatrix} =: J_f(\boldsymbol{x}) =: \frac{\partial f}{\partial \boldsymbol{x}}$$

Multivariate Taylor series: Given differentiable f, it holds that

$$f(\boldsymbol{x}) = f(\boldsymbol{x}_0) + J_f(\boldsymbol{x}_0)(\boldsymbol{x} - \boldsymbol{x}_0) + o\big(\|\boldsymbol{x} - \boldsymbol{x}_0\|\big)$$

Multivariate chain rule: Given differentiable $f : \mathbb{R}^n \to \mathbb{R}^m$ and $g : \mathbb{R}^m \to \mathbb{R}^p$, it holds that

$$\frac{d}{d\boldsymbol{x}} g\big(f(\boldsymbol{x})\big) = J_{g \circ f}(\boldsymbol{x}) = J_g\big(f(\boldsymbol{x})\big) J_f(\boldsymbol{x}) = \frac{\partial g}{\partial f} \frac{\partial f}{\partial \boldsymbol{x}}$$

**Examples**

$$z = z(x) = Wx \qquad \Rightarrow \quad J_z(x) = W$$

$\overbrace{\phantom{J_z(x)}}^{\frac{\partial z}{\partial x}}$

$$\varphi(x) = \left[\, \varphi(x_1) \, , \cdots , \, \varphi(x_n) \right]^T \quad \Rightarrow \quad J_\varphi(x) = \overbrace{\begin{pmatrix} \varphi'(x_1) & 0 & \text{---} & 0 \\ & \varphi'(x_2) & & 0 \\ 0 & & \ddots & \varphi'(x_n) \end{pmatrix}}^{diag(\varphi'(x))}$$

$\varphi'(x) = \left[\varphi'(x_1) \cdots \varphi'(x_n)\right]$

$$f(x) = \varphi(\underbrace{Wx}_{z}) \quad \Rightarrow \quad J_f(x) = J_\varphi(z) \, J_z(x) = diag(\varphi'(z)) \cdot W$$

# Backpropagation

Consider neural network $f(x; W)$ with weights $W = [W^{(1)}, W^{(2)}, ..., W^{(L)}]$

Want to compute gradient of loss $\ell(W; x, y)$ w.r.t., weights $W^{(i)}$

$$\left(\nabla_{W^{(L)}} \ell\right)^T = \frac{\partial \ell}{\partial W^{(L)}} = \frac{\partial \ell}{\partial f} \frac{\partial f}{\partial W^{(L)}}$$

$$\left(\nabla_{W^{(L-1)}} \ell\right)^T = \frac{\partial \ell}{\partial W^{(L-1)}} = \frac{\partial \ell}{\partial f} \frac{\partial f}{\partial z^{(L-1)}} \frac{\partial z^{(L-1)}}{\partial W^{(L-1)}}$$

$$\left(\nabla_{W^{(L-2)}} \ell\right)^T = \frac{\partial \ell}{\partial W^{(L-2)}} = \frac{\partial \ell}{\partial f} \frac{\partial f}{\partial z^{(L-1)}} \frac{\partial z^{(L-1)}}{\partial z^{(L-2)}} \frac{\partial z^{(L-2)}}{\partial W^{(L-2)}}$$

$$\vdots$$

$$\left(\nabla_{W^{(i)}} \ell\right)^T = \frac{\partial \ell}{\partial W^{(i)}} = \frac{\partial \ell}{\partial f} \frac{\partial f}{\partial z^{(L-1)}} \frac{\partial z^{(L-1)}}{\partial z^{(L-2)}} \cdots \frac{\partial z^{(i+1)}}{\partial z^{(i)}} \frac{\partial z^{(i)}}{\partial W^{(i)}}$$

Key insight: Can reuse computations from forward propagation and from layer i+1 to compute $W^{(i)}$