# Lecture 10

Backprop

Sol 1,2,3

Q1 : 4:00 Thu 2/12
        15 min

L$\overline{\mathbb{N}}$ - 2

# Lecture Notes IV – Neural Networks, Part 2

Marina Meilă

mmp@uwaterloo.ca

February 8, 2026

# Backpropagation

Train n.n. by G.D.

Training a single unit ←

Training a 2-layer network ←

Training a *L*-layer network ←

**Reading** HTF Ch.: 11.3 Neural networks, Murphy Ch.: (16.5 neural nets), Bach Ch.: –, Deep Learning Book (Goodfellow, Bengio, Courville) 6.1-4, ResNet 7.6, ConvNet 9., Autoencoders 14.1, Dive Into Deep Learning 4.1-4.3.

## Training a neural network    **BIG PICTURE**

- Model is multilayer network.
- Layers $l = 1, 2, \ldots L$, with $x^{(l)} \in \mathbb{R}^{m_l}$, for $l = 1 : L-1$, $m_0 = d$, $m_L = 1$ (for regression and binary classification).

$$
\begin{align}
x^{(0)} &= x \tag{1}\\
x^{(L)} &= f(x) \tag{2}\\
x^{(l)} &= \phi(z^{(l)}) \quad \text{for } l = 1 : L-1 \tag{3}\\
x^{(L)} &= z^{(L)} \equiv \phi_{\text{out}}(z^{(L)1}) \tag{4}\\
z^{(l)} &= W^{(l)} x^{(l-1)} \tag{5}\\
&\tag{6}
\end{align}
$$

- Parameters $\mathbb{W} = \{W^{(l)} \in \mathbb{R}^{m_l \times m_{l=1}}\}$
- $W_k^{(l)}$ is row $k$ of $W^{(l)}$ and corresponds to unit $k$ of layer $l$
- $\mathcal{D} = \{(x^1, y^1), \ldots (x^n, y^n)\}$

$$W_i^{(\ell)} \sim N\left(0, \frac{v^2}{m_\ell} I\right)$$
$$\uparrow$$
$$\mathbb{R}^{m_\ell}$$

- **How to train this network?**
- Minimize $\mathcal{L}(\mathbb{W})$ (remember $\phi_{\text{out}}$ is associated with $\mathcal{L}$)
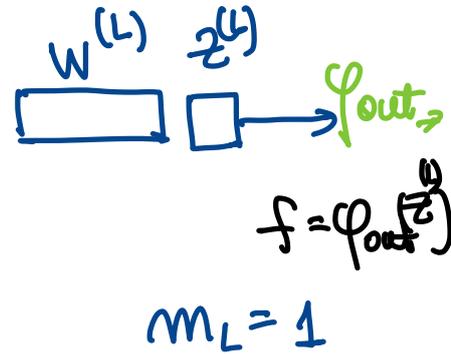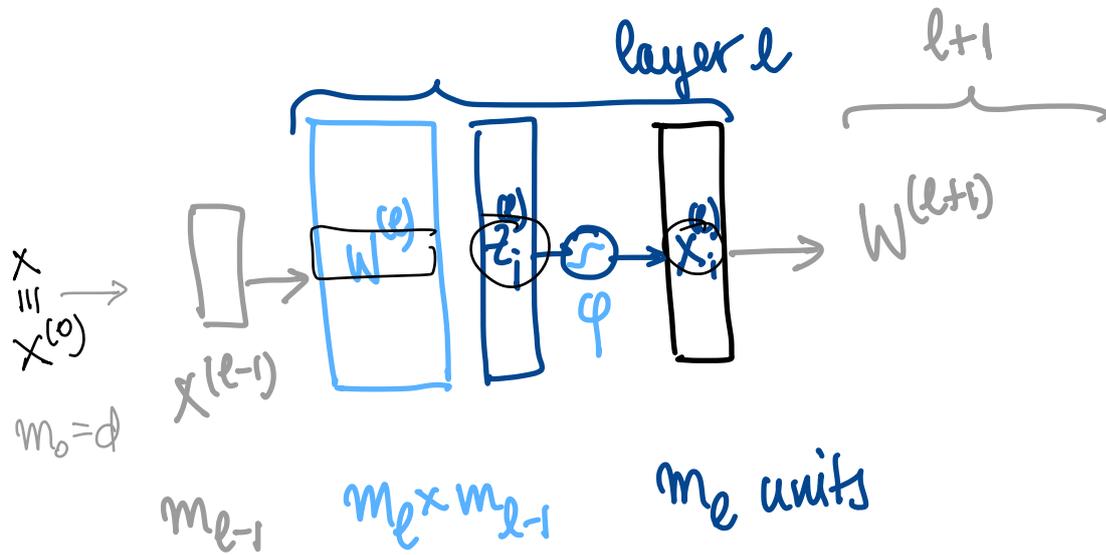
- Minimization by Gradient Descent
  - Initialize $\mathbb{W}$ to small random values **1.**
  - for $t = 1, 2, \ldots$ do $\mathbb{W} \leftarrow \mathbb{W}^t - \eta \frac{\partial \mathcal{L}}{\partial \mathbb{W}} (\mathbb{W}^t)$ **2.**

- **We need to compute gradient** $\frac{\partial \mathcal{L}(\mathbb{W})}{\partial w_{ij}^{(l)}}$ for all $i = 1 : m_l$, $j = 1 : m_{l-1}$, $l = 1 : L$.

$$x^{(0)} = x$$

$$m_0 = d$$

layer $\ell$

$\ell + 1$

$W^{(\ell+1)}$

$x^{(\ell-1)}$

$W^{(\ell)}$

$z_i^{(\ell)}$

$\varphi$

$x_i^{(\ell)}$

$m_{\ell-1}$

$m_\ell \times m_{\ell-1}$

$m_\ell$ units

$W^{(L)}$  $z^{(L)}$

$\varphi_{out}$

$f = \varphi_{out}[z^{(L)}]$
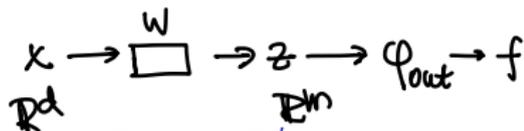
$m_L = 1$

$\ell = 1 : L$

$$z^{(\ell)} = W^{(\ell)} x^{(\ell-1)}$$

$$\varphi(z_i^{(\ell)}) = x_i^{(\ell)}$$

$$\left[ \text{parameters } W = \{ W^{(1)}, W^{(2)}, \ldots W^{(L)} \} \right]$$

# Training a single unit   *single* $(x, y)$

$$x \rightarrow \boxed{\ \ }^{w} \rightarrow z \rightarrow \phi_{out} \rightarrow f$$
$$\mathbb{R}^d \qquad \mathbb{R}^m$$

▶ The function $f(x; \mathbb{W}) = \phi_{out} \underbrace{\left( \sum_{j=1}^{d} w_j x_j \right)}_{z}$ with parameters $\mathbb{W} = w \in \mathbb{R}^d$.

▶ The loss function $\mathcal{L}(y, f(x; \mathbb{W})) =$ Least Squares, or Logistic (Max Likelihood) binary or multiclass.
   **Fact** For each $\mathcal{L}$, there exists a mapping $y \rightarrow y_*$, $z \rightarrow \phi_{out}(z)$ such that

   $$\boxed{-\frac{\partial \mathcal{L}(y, f(x; \mathbb{W}))}{\partial f} = y_* - \phi_{out}} = \text{target output} - \text{network output}$$

▶ Let's use  chain rule

$$y \rightarrow y_*$$

$$z = wx$$

$$L = \frac{1}{2}\left( y_* - \phi_{out} \right)^2$$
$$\text{LS}$$

$$\frac{\partial L}{\partial \phi_{out}} = \phi_{out} - y_*$$

$$-\frac{\partial \mathcal{L}}{\partial f} = y_* - \phi_{out}(z) \tag{7}$$

$$\frac{\partial f}{\partial z} = \phi'_{out}(z) \tag{8}$$

$$\frac{\partial z}{\partial w}$$

$$\frac{\partial f}{\partial w} = \phi'_{out}(z)x \tag{9}$$

$$\frac{\partial f}{\partial x} = \phi'_{out}(z)w \tag{10}$$

$$\frac{\partial z}{\partial x}$$

$$-\frac{\partial \mathcal{L}}{\partial w} = (y_* - \phi_{out}(z))\phi'_{out}(z)x$$

▶ Remember $\phi'_{out}(z) = z$ for $\mathcal{L}_{LS}$, $\phi'_{out}(z) = \phi(x)(1 - \phi(z))$ for $\mathcal{L}_{logit}$.
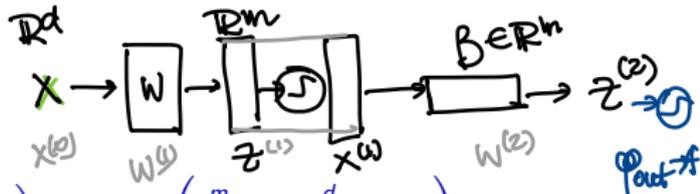▶ Notice that $\partial \mathcal{L}/\partial w$ is a vector collinear with $x$.

# Training a single unit, *n* data points

- $\mathcal{D} = \{(x^1, y^1), \ldots (x^n, y^n)\}$
- For single pair $(x, y)$, $\frac{\partial L}{\partial w} = (y - \phi_{\text{out}}(z))\phi'_{\text{out}}(z)x$
- For entire $\mathcal{D}$, $\mathcal{L}(\mathbb{W}) = \frac{1}{n}\sum_{i=1}^{n}\mathcal{L}(y^i, f(x^i; \mathbb{W}))$.
- The gradient of this loss is

$$-\frac{\partial L}{\partial w} = \frac{1}{n}\sum_{i=1}^{n}\left(\frac{\partial}{\partial w}\mathcal{L}(y^i, f(x^i; \mathbb{W}))\right) = \frac{1}{n}\sum_{i=1}^{n}(y^i - \phi_{\text{out}}(z^i))\phi'_{\text{out}}(z^i)x^i. \quad (11)$$

## Training a 2-layer network

► Consider a two layer neural network

$$f(x) = \phi_{\text{out}}\left(\sum_{i=1}^{m} \beta_i x_i^{(1)}\right) = \phi_{\text{out}}\left(\sum_{i=1}^{m} \beta_i \phi(\sum_{j=1}^{d} w_{ij} x_j)\right) \tag{12}$$

The parameters $\mathbb{W}$ are $\beta$ and $W = [w_{ij}]_{i=1:m, j=1:d}$

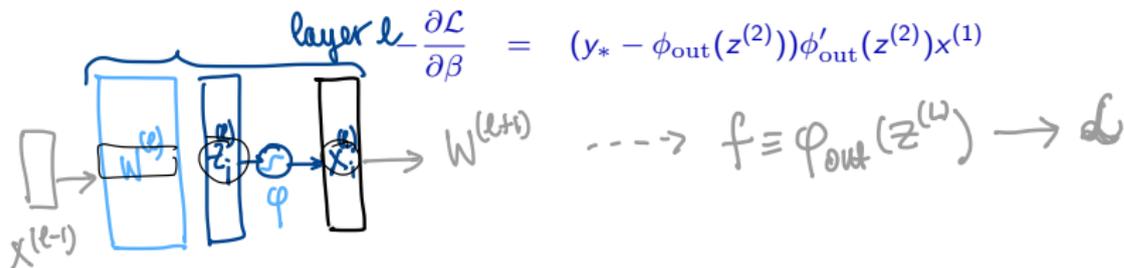**Gradient w.r.t.** $\beta$ (same as before! $w \leftarrow \beta$)

$$-\frac{\partial \mathcal{L}}{\partial f} = y_* - \phi_{\text{out}}(z^{(2)}) \tag{13}$$

$$\frac{\partial f}{\partial x^{(1)}} = \phi'_{\text{out}}(z^{(2)})\beta \tag{14}$$

$$\frac{\partial f}{\partial x_i^{(1)}} = \phi'_{\text{out}}(z^{(2)})\beta_i \tag{15}$$

$$\frac{\partial f}{\partial \beta} = \phi'_{\text{out}}(z^{(2)})x^{(1)} \tag{16}$$

$$-\frac{\partial \mathcal{L}}{\partial \beta} = (y_* - \phi_{\text{out}}(z^{(2)}))\phi'_{\text{out}}(z^{(2)})x^{(1)}$$

# Trainining a 2-layer network – hidden layer parameters

**Gradient w.r.t. $W$**

- Let's break $W \in \mathbb{R}^{m \times d}$ into rows. $W_i = [w_{ij}]_{j=1:d}$ is the column vector of weights for unit $i$ in hidden layer.

- We have $z_i = W_i^T x$ and $x_i^{(1)} = \phi(z_i)$.

- Again, as before, with $\phi_{\text{out}} \leftarrow \phi$, $w \leftarrow W_i$, $z \leftarrow z_i$, $z^{(2)} = \beta^T x_i^{(1)}$, $f = \phi_{\text{out}}(z^{(2)})$ we have

$$\frac{\partial x_i^{(1)}}{\partial z_i} = \phi'(z_i) \qquad (17)$$

$z_i = W_i x$

$$\frac{\partial x_i^{(1)}}{\partial W_i} = \phi'(z_i)x \qquad (18)$$

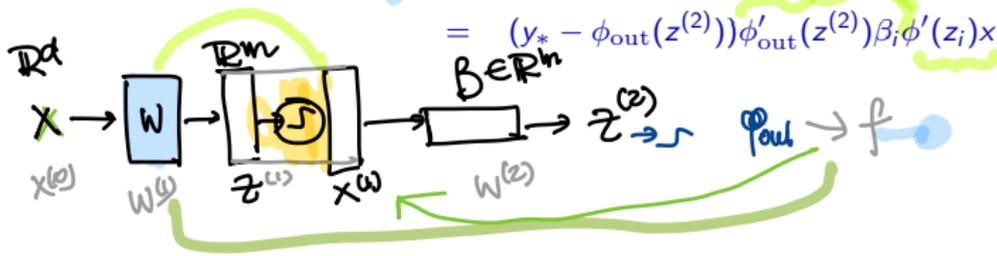$\frac{\partial z_i}{\partial x} = W_i$

$$\frac{\partial f}{\partial W_i} = \frac{\partial f}{\partial x_i^{(1)}} \frac{\partial x_i^{(1)}}{\partial W_i} \qquad (19)$$

$$= \phi'_{\text{out}}(z^{(2)})\beta_i \phi'(z_i)x \qquad (20)$$

$$-\frac{\partial \mathcal{L}}{\partial W_i} = \frac{\partial \mathcal{L}}{\partial f} \frac{\partial f}{\partial W_i} \qquad (21)$$

$$= (y_* - \phi_{\text{out}}(z^{(2)}))\phi'_{\text{out}}(z^{(2)})\beta_i \phi'(z_i)x \qquad (22)$$

$$\qquad (23)$$

# From 2 layers to $L$ layers

$|z|$ large $\Rightarrow \phi' \approx 0$ "saturation"



$$\frac{\partial x_i^{(l)}}{\partial z_i^{(l)}} = \phi'(z_i^{(l)}) \tag{24}$$

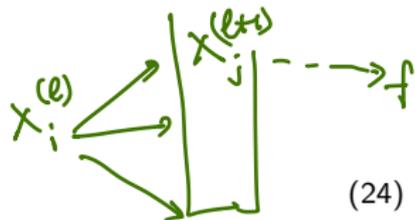$$\frac{\partial x_i^{(l)}}{\partial W_i^{(l)}} = \phi'(z_i^{(l)}) x^{(l-1)} \tag{25}$$

$$z_i^{(\ell)} = W_i^{(\ell)} x^{(\ell-1)}$$

used for $\left\{ \frac{\partial x_i^{(l)}}{\partial x_i^{(l-1)}} = \phi'(z_i^{(l)}) W_i^{(l)} \right. \tag{26}$

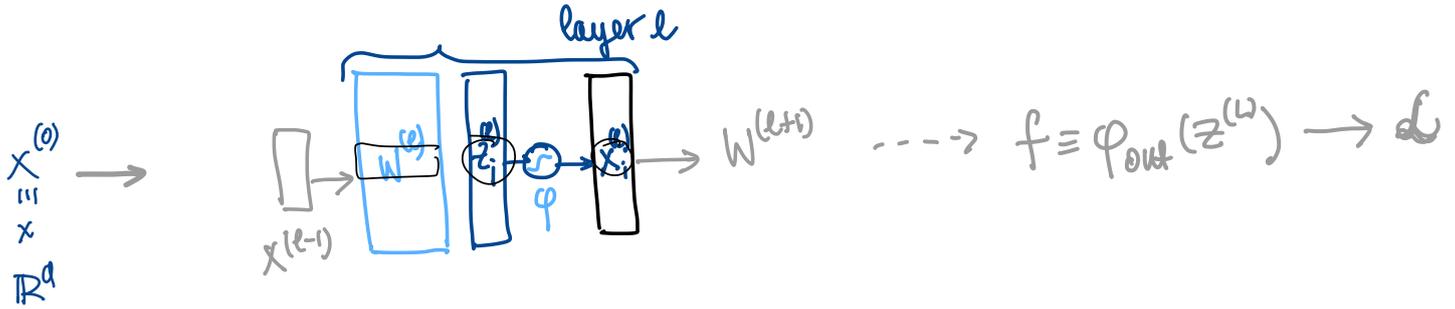$\frac{\partial \mathcal{L}}{\partial W^{(\ell-1)}}$

$$\frac{\partial f}{\partial x_i^{(l)}} = \sum_{j=1}^{n_{l+1}} \frac{\partial f}{\partial x^{(l+1)}} \frac{\partial x_j^{(l+1)}}{\partial x_i^{(l)}} \tag{27}$$

$$\frac{\partial f}{\partial W_i^{(l)}} = \frac{\partial f}{\partial x_i^{(l)}} \frac{\partial x_i^{(l)}}{\partial W_i^{(l)}} \tag{28}$$

layer $\ell$ $\quad \frac{\partial \mathcal{L}}{\partial W_i^{(l)}} = (y_* - \phi_{\text{out}}) \frac{\partial f}{\partial W_i^{(l)}}$



(29)

$x^{(l-1)}$ $\qquad W^{(\ell+i)} \quad \dashrightarrow \quad f \equiv \phi_{\text{out}}(z^{(L)}) \longrightarrow \mathcal{L}$

$x^{(0)}$
$\equiv$
$x$
$\mathbb{R}^q$



layer $\ell$

$W^{(\ell)}$  $z_i^{(\ell)}$  $\varphi$  $x_i^{(\ell)}$  $\longrightarrow$ $W^{(\ell+i)}$  $- - - \to$  $f \equiv \varphi_{out}(z^{(L)}) \longrightarrow \mathcal{L}$

$x^{(\ell-1)}$

for $t = 1, 2, \cdots$

    for $k = 1:n$ (examples)

FORWARD

(Prediction

     for $\ell = 1 : L$    $\downarrow$ store!

       compute $x^{(\ell)}$, $z^{(\ell)}$ from $x^{(\ell-1)}$, $W^{(\ell)}$

       $f = \varphi_{out}(z^{(L)})$, $\mathcal{L}^{train} += \frac{1}{n} \mathcal{L}(y^i, f(x^i))$ $\} \to$ used for stopping

                                                      not $\frac{\partial \mathcal{L}}{\partial W}$

                                    BACK Prop.

     for $\ell = L : -1 : 1$   $\frac{\partial \mathcal{L}}{\partial ...}$

       compute $\frac{\partial \mathcal{L}}{\partial W_i^{(\ell)}}$ from $x^{(\ell)}, \frac{\partial f}{\partial x^{(\ell)}}, W^{(\ell)}, \frac{\partial x^{(\ell+1)}}{\partial x^{(\ell)}}$   $\frac{\partial \mathcal{L}}{\partial f}$

       $i = 1 : m_\ell$

            $\frac{\partial x^{(\ell)}}{\partial x^{(\ell-1)}} \to$ for $\frac{\partial \mathcal{L}}{\partial W}^{(\ell-1)}$