# Lecture 12

Linear regression

$$X_1 X_2 \ldots X_d \longrightarrow Y$$

$\nwarrow$ real-valued

Logistic regression

$$X_1 X_2 \ldots X_d \longrightarrow Y$$

$\nwarrow$ categorical or symbolic

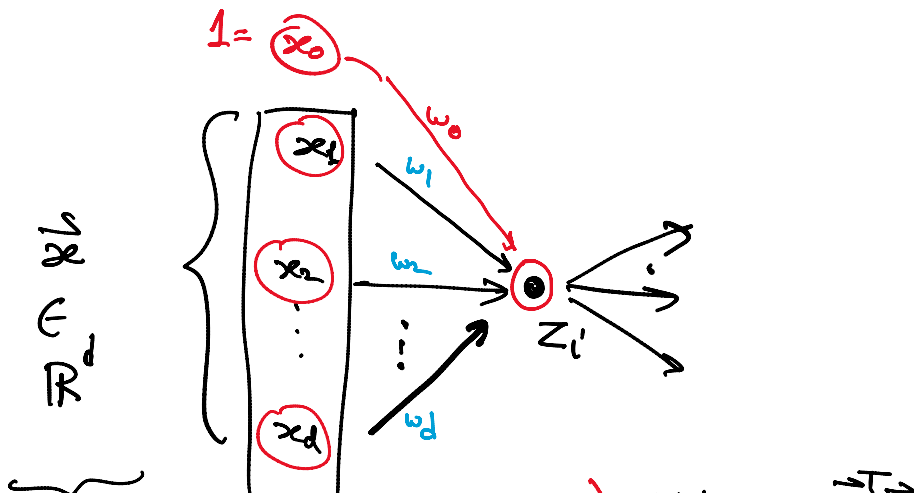$$\{c_1, c_2, \ldots, c_k\}$$

linear/identity

$$f(X_1, X_2, X_d) = \omega_1 X_1 + \omega_2 X_2 + \cdots + \omega_d X_d + \omega_0 = \underline{\underline{\tilde{\omega}^T \tilde{X}}}$$

$$f(X_1, X_2, X_d) = \boxed{\text{softmax}} \left( \tilde{\omega}^T \tilde{X} \right) \quad \text{or} \quad \boxed{\text{Sigmoid}} \left( \tilde{\omega}^T \tilde{X} \right)$$

non-linear        linear        binary case

---

## MLP: Multilayer Perceptrons  ( Feed Forward Neural Nets)

### Multi-layer logistic regression        (artificial NN)



$$1 = \textcircled{$\tilde{x}_0$}$$

$$\omega_0$$

$x_1$    $\omega_1$

$x_2$    $\omega_2$

$x \in \mathbb{R}^d$

$x_d$    $\omega_d$

$z_i$

$\tilde{\omega}^T \tilde{x}$

$$\underbrace{\qquad}_{\substack{\text{input} \\ \text{data points}}} \quad \fbox{$\boxed{x_d}$} \; {}^{/w_d}$$

1) $\text{Net}_{z_i} = \vec{w}^T \vec{x} = w_0 + w_1 x_1 + \dots + w_d x_d$

2) activation function

$$z_i = f(\text{Net}_{z_i})$$

0) 
Step function

$$z_i = f(net_i) = \begin{cases} 1 & \text{if } \boxed{net_i > 0} \\ 0 & \text{otherwise} \end{cases}$$

$w_1 x_1 + w_2 x_2 + \dots + w_d x_d > -w_0$

$\underbrace{\qquad}_{\text{threshold}}$

Common Activation functions

1) Linear / Identity

$$f(net_i) = net_i$$

2) Sigmoid

$$f(net_i) = \frac{1}{1 + e^{-net_i}}$$

output $[0, 1]$

3) tanh

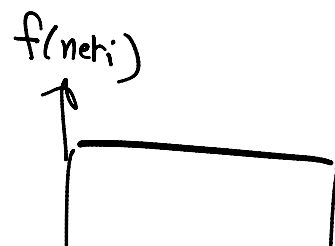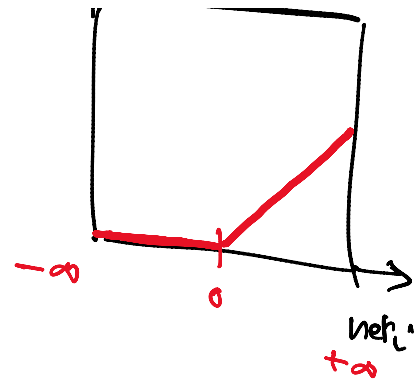$$f(net_i) = \frac{e^{net_i} - e^{-net_i}}{e^{net_i} + e^{-net_i}}$$

output $[-1, 1]$

4) ReLU (rectified linear)

$f(\dots)$

$\begin{cases} net_i & \text{if } net_i > 0 \end{cases}$

$f(net_i)$

$$f(net_i) = \begin{cases} net_i & \text{if } net_i > 0 \\ 0 & \text{otherwise} \end{cases}$$



$$= \max\{0, net_i\}$$

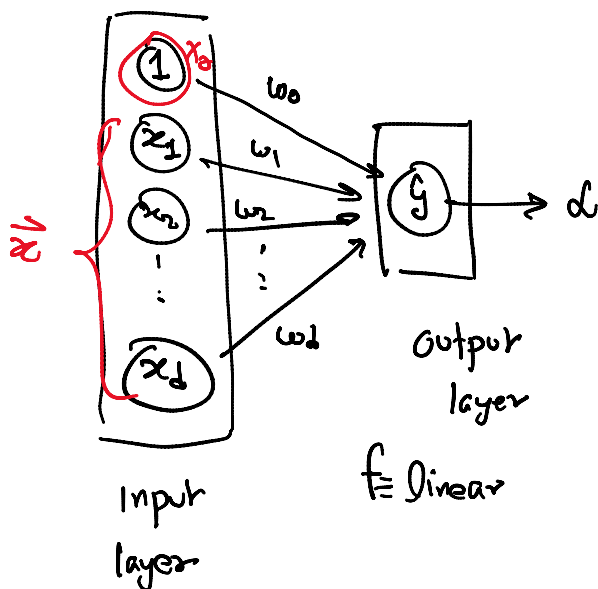5) $softmax\left(net_i \mid \underbrace{net_1, net_2, \ldots, net_m}_{m \text{ different neurons}}\right)$

<span style="color:red">Output prob distribuhón over the m neurons!</span>

$$= \frac{e^{net_i}}{\sum\limits_{j=1}^{m} e^{net_j}}$$

## Linear Network (Linear regression)



Input layer

Output layer

$f \equiv$ linear

$$net_{\hat{y}} = \tilde{\omega}^T \tilde{x}$$
$$= \omega_0 + \omega_1 x_1 + \ldots + \omega_j x_d$$

$$\hat{y} = linear(net_i) = net_i$$

$\tilde{\omega}$ is unknown!

$f \equiv$ linear

learn: how?

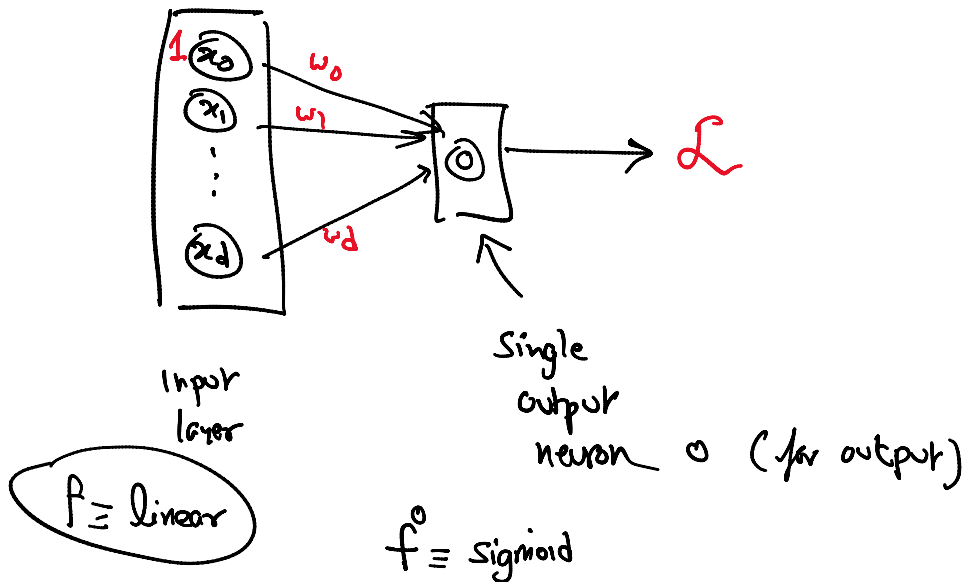$$\mathcal{L} = \text{squared error} = \left(y - \hat{y}\right)^2$$

true ↗ ↖ output of NN

$$\nabla_{\tilde{\omega}} = \frac{\partial \mathcal{L}}{\partial \tilde{\omega}}$$

train over a mini-batch of data via gradient descent.

---

Logistic regression network



Input layer

$f \equiv$ linear

Single output neuron $o$ (for output)

$f^{o} \equiv$ sigmoid

$$net_o \leftarrow \omega_0 x_0 + \omega_1 x_1 + \cdots + \omega_d x_d$$
$$= \tilde{\omega}^T \tilde{x}$$

$$o = f^{o}(net_o) = sigmoid(net_o)$$

$$\mathcal{L} \equiv \text{binary cross entropy loss}$$
$$= -\left(y \log(o) + (1-y) \log(1-o)\right)$$

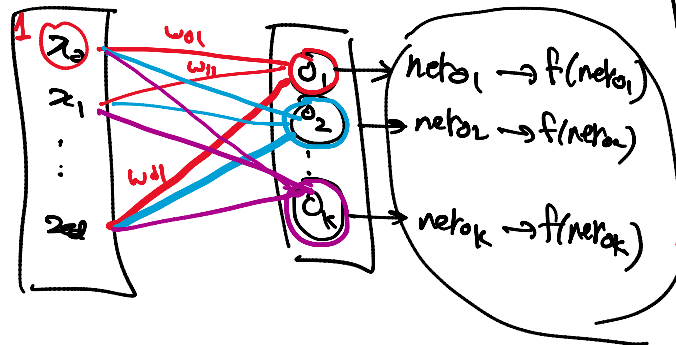$$y \in \{0, 1\} \quad \text{true class}$$

$$\nabla_{\tilde{w}} = \frac{\partial \mathcal{L}}{\partial \tilde{w}} \qquad\qquad 0 \in [0,1]$$

---



K-class logistic regression

$w_{ij} \equiv$ weight between input $x_i$ and output $o_j$

Input layer          output layer

$$W = \begin{bmatrix} \begin{matrix} w_{01} \\ w_{11} \\ \vdots \\ w_{d1} \end{matrix} & \begin{matrix} w_{02} \\ w_{12} \\ \vdots \\ w_{d2} \end{matrix} & \cdots & \begin{matrix} w_{0k} \\ w_{1k} \\ \vdots \\ w_{dk} \end{matrix} \end{bmatrix}$$

$d \times k$ matrix

$\overset{o}{f} \equiv$ softmax

Output of the network

$$\begin{cases} O_1 = \overset{o}{f}(\,net_1\, | \, net_1 \; net_2 \cdots net_k\,) \\ \qquad\qquad \uparrow \\ \qquad\quad softmax \\ O_2 = softmax(net_2) \\ \quad\vdots \\ O_k = softmax(net_k) \end{cases}$$

these are the probabilities for each class

$$\mathcal{L} = f(O_1, O_2 \cdots O_k, \vec{y})$$

One-hot class vector

$$\mathcal{L} \equiv CE \equiv \text{Cross entropy}$$

$$= -\sum_{i}^{k} y_i \log(O_i)$$

$$= - \sum_{i=1} y_i \log(o_i)$$

$$\nabla_W = \frac{\partial \mathcal{L}}{\partial W}$$

SGD ← Stochastic "mini-batch" GD



$$\vec{b} = \begin{pmatrix} w_{01} \\ w_{02} \\ \vdots \\ w_{0d} \end{pmatrix}$$

$$W = \{w_{ij}\}_{\substack{i=1\dots d \\ j=1\dots k}}$$

$d \times k$

Two parameters $\Theta = \{W, \vec{b}\}$

learn via backprop (SGD)

---

Simplest extension:

add one <u>hidden</u> layer



input      hidden      output

linear : regression ] Value

sigmoid : binary logistic

softmax : k-way logistic ] prob

input          hidden          output

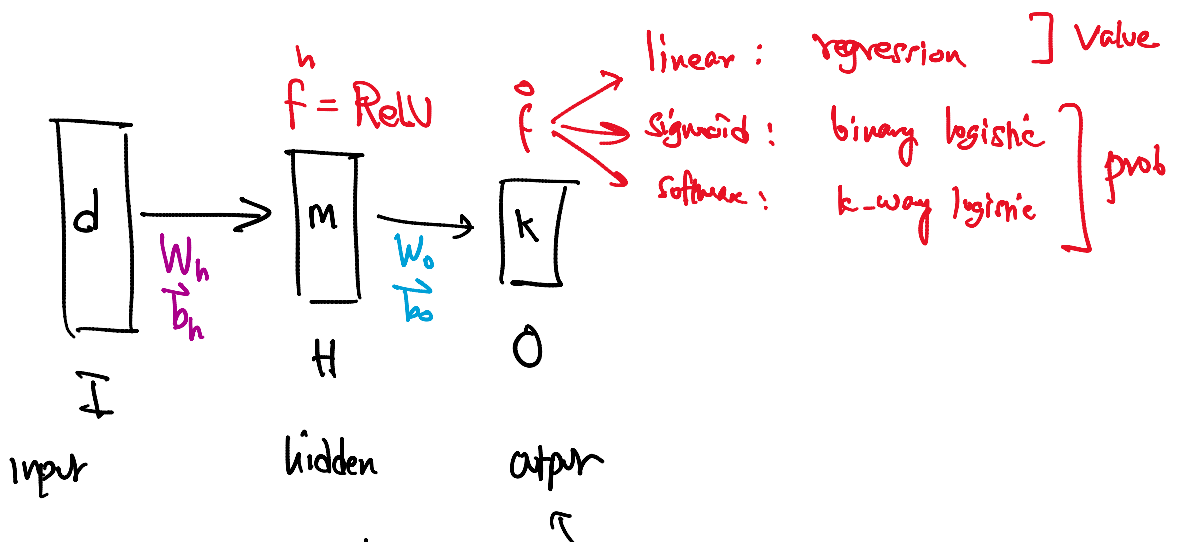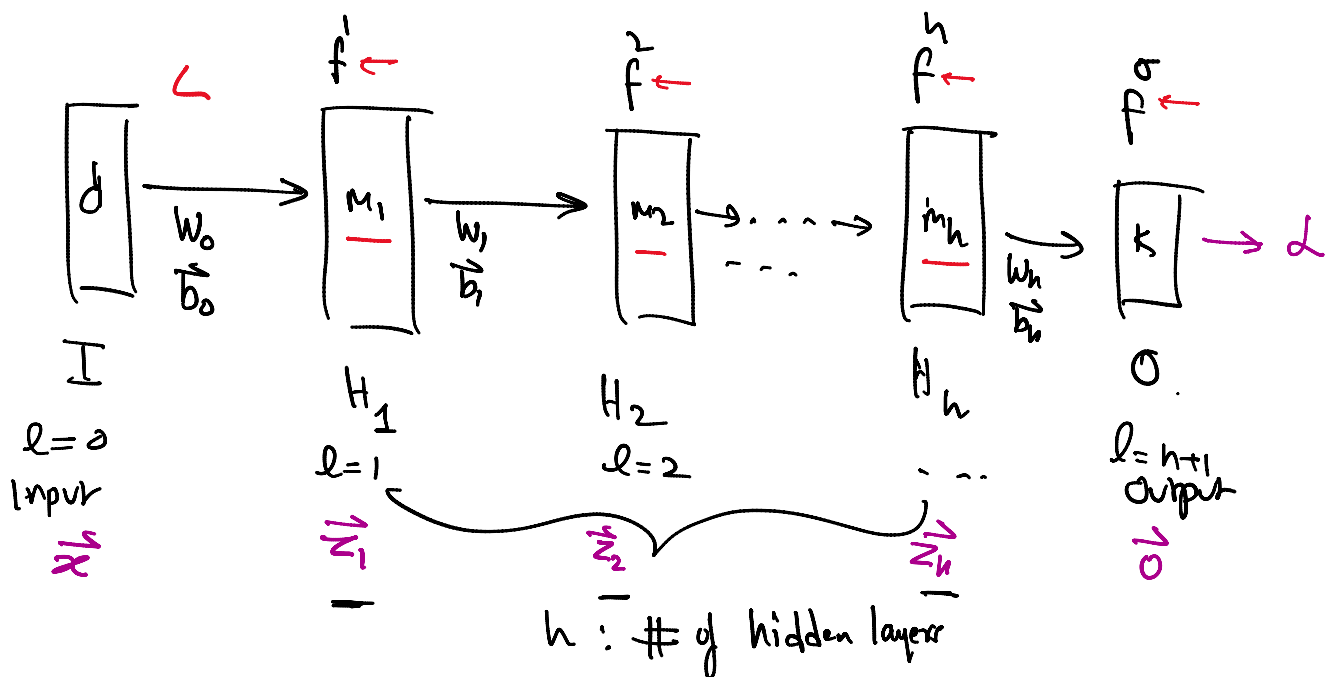$d$ : # of          M: we have          k: determine by
   independent          to choose          nature of Y
   vars

$(X_1, X_2 \cdots X_d)$

between I and H
$W_h$ : $d \times m$     matrix of synaptic weights

$\vec{b}_h$ : $m$          bias vector for H

between H & O
$W_o$ : $M \times K$

$\vec{b}_o$ : $k$

← bias for O

## MLP



$f^1$          $f^2$          $f^h$          $f^\sigma$

$d$  →  $M_1$  →  $M_2$  ⇢ $\cdots$ ⇢  $m_h$  →  $k$  → $\alpha$
   $W_o$          $W_1$                    $W_h$
   $\vec{b}_0$     $\vec{b}_1$              $\vec{b}_h$

I          $H_1$          $H_2$          $H_h$          O

$\ell = 0$      $\ell = 1$      $\ell = 2$          $\ell = h+1$
Input                                               output

$\vec{x}$      $\vec{z}_1$      $\vec{z}_2$      $\vec{z}_h$      $\vec{o}$

$h$ : # of hidden layers

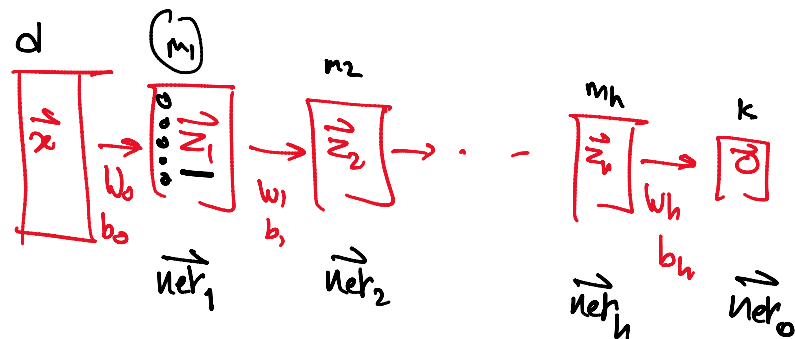$W_i$ is the weight matrix between layer $i$  and $i+1$

$\vec{b}_i$ is the bias for $i+1$

$$\Theta = \{ W_i, \vec{b}_i \mid i = 0 \cdots h \}$$

use back propagation

① Feed - forward step

$\vec{x} \longrightarrow$  how do I compute  $\vec{o}$



h: # of hidden layers.

input : $\vec{x}$

$H_1 \begin{cases} \vec{net_1} = W_0^T \vec{x} + \vec{b_0} \\[2mm] \vec{z_1} = f^1(\vec{net_1}) = f\left( \underbrace{\underbrace{W_0^T}_{m_1 \times d} \underbrace{\vec{x}}_{d \times 1}}_{M_1 \times 1} + \underbrace{\vec{b_0}}_{m_1 \times 1} \right) \end{cases}$

$W_0 : d \times m_1$
$\vec{b_0} : m_1$

$M_1 \times 1$

$H_2 \begin{cases} \vec{net_2} = W_1^T \vec{z_1} + \vec{b_1} \\[2mm] \vec{z_2} = f^2(\vec{net_2}) \end{cases}$

$\vdots$

$H_h \begin{cases} \vec{net_h} = W_{h-1}^T \vec{z_{h-1}} + \vec{b_{h-1}} \\[2mm] \vec{z_h} = f^h(\vec{net_h}) = f^h\left( W_{h-1}^T \vec{z_{h-1}} + \vec{b_{h-1}} \right) \end{cases}$

$$\cdots _n - T(net_h) = F(W_{h-1} \vec{z}_{h-1} + b_{h-1})$$

$$0 \begin{cases} \vec{net_0} = W_h^T \vec{z}_h + \vec{b}_h \\ \vec{0} = \mathring{f}(\vec{net_0}) = \mathring{f}(W_h^T \vec{z}_h + \vec{b}_h) \end{cases}$$

$$\mathcal{E} = \mathcal{L} \longrightarrow \text{squared error} = \frac{1}{2} \|\vec{y} - \vec{o}\|^2 = \mathcal{L}$$

error ↑   loss ↑   ↘ cross entropy ↗ binary

↘ k-way

$$\mathcal{L} = -\sum_{i=1}^{k} y_i \log(o_i)$$

$$\vec{0} = (o_1, o_2, \ldots o_k)^T$$

$$\vec{y} = (y_1, y_2, \ldots, y_k)^T$$

⎵ One-hot vector

Initialization: randomly initialize all

$$W_i \quad \text{and} \quad \vec{b}_i \qquad i = 0 \ldots h$$

(small, random normal)

---
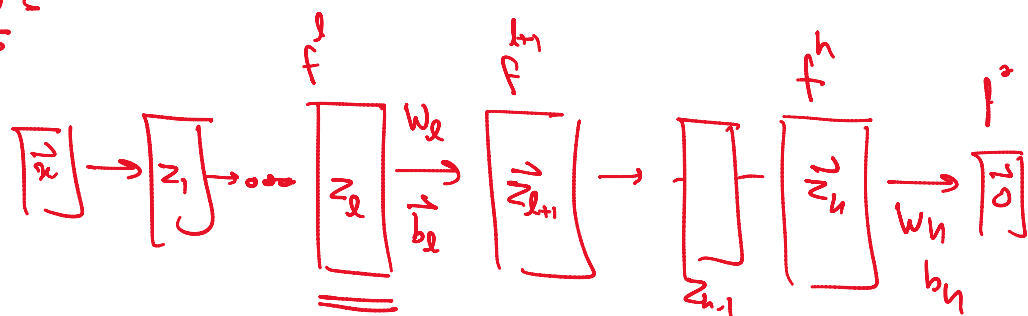
2) Backpropagation step

(update the $W_i, \vec{b}_i$ based on $\mathcal{L}$)

what: $\boxed{\nabla_{W_i} = \dfrac{\partial \mathcal{L}}{\partial W_i}} \qquad \boxed{\nabla_{\vec{b}_i} = \dfrac{\partial \mathcal{L}}{\partial \vec{b}_i}} \quad \forall i = 0 \ldots h$

$$W_i^{(t+1)} = W_i^{(t)} - \eta \cdot \nabla_{W_i} \qquad \vec{b}_i^{(t+1)} = \vec{b}_i^{(t)} - \eta \cdot \nabla_{b_i}$$

How 2

$$\vec{x} \rightarrow \boxed{z_1} \rightarrow \cdots \quad \overset{f^{\ell}}{\boxed{z_{\ell}}} \overset{W_{\ell}}{\underset{\vec{b}_{\ell}}{\rightarrow}} \overset{f^{h_{-1}}}{\boxed{\vec{z}_{\ell+1}}} \rightarrow \quad \boxed{\underset{z_{h-1}}{\phantom{x}}} \quad \overset{f^{h}}{\boxed{\vec{z}_{h}}} \underset{\substack{W_h \\ b_h}}{\rightarrow} \overset{f^{o}}{\boxed{\vec{o}}}$$

$$z_{\ell} = f^{\ell}\left(W_{\ell-1}^{T} \vec{z}_{\ell-1} + \vec{b}_{\ell-1}\right)$$

$$\vdots$$

$$\vec{o} = f^{o}\left(W_{h}^{T} \;\boxed{\vec{z}_{h}}\; + \vec{b}_{h}\right)$$

$$= f^{o}\left(W_{h}^{T} f^{h}\left(W_{h-1}^{T}\vec{z}_{h-1} + \vec{b}_{h-1}\right) + \vec{b}_{h}\right)$$

$$\ell = \tfrac{1}{2}\|\vec{y} - \vec{o}\|^2$$

$$\frac{\partial \ell}{\partial W_{\ell}} = \tfrac{1}{2}\left\|\vec{y} - \left(\left(\left(\phantom{xxx}\right)\right)\right)\right\|^{2}$$

$$\underbrace{\phantom{xxxxxxxxxxx}}_{\text{chain-rule}}$$