# Deep Learning
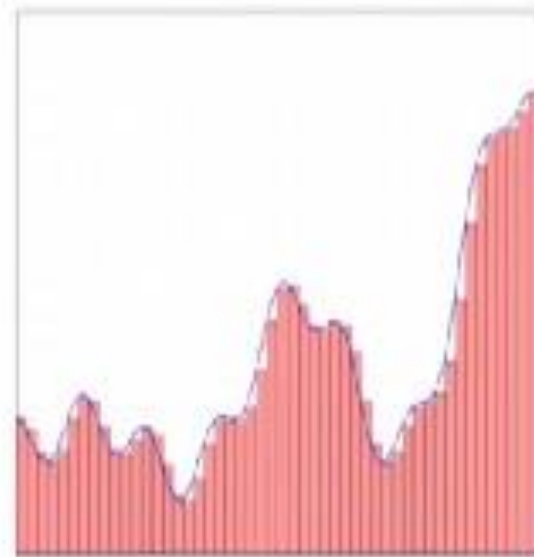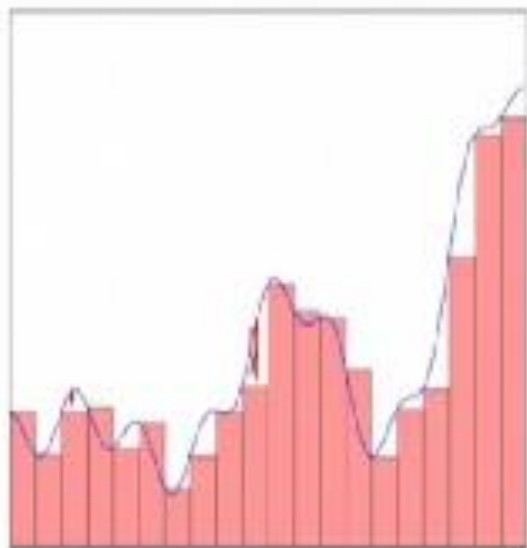
# Representative Power of Multilayer Networks

- A multilayer network of perceptrons with a single hidden layer can be used to approximate any Boolean function precisely

- A multilayer network of sigmoid neurons with a single hidden layer can be used to approximate any continuous function to any desired precision

Sigmoid

# Multilayer Network

- For any function $f(x): \mathbb{R}^n \longrightarrow \mathbb{R}^m$, we can find a network with enough neurons, whose output $g(x)$ satisfies $|g(x) - f(x)| < \epsilon$
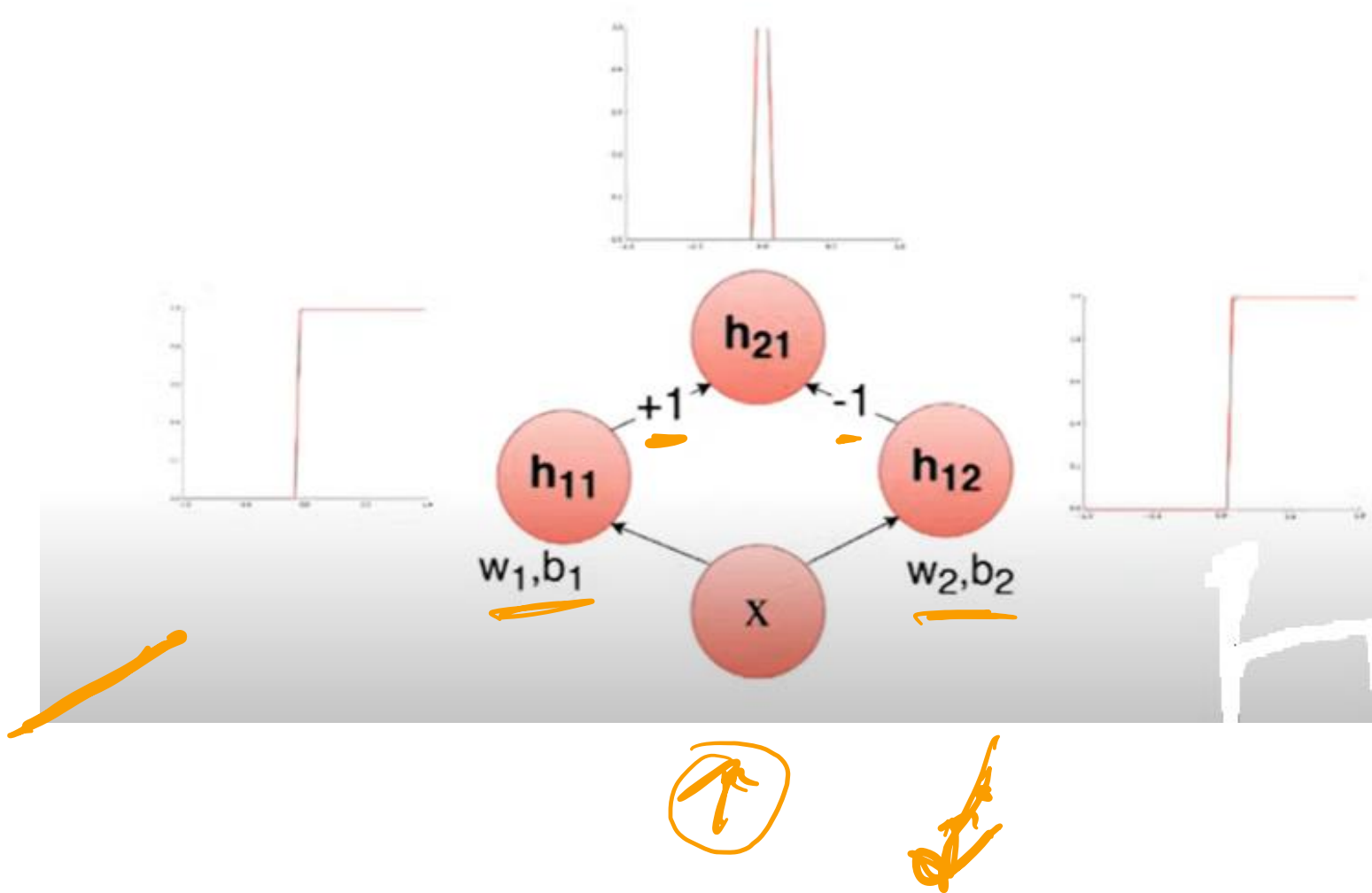- Such an arbitrary function can be represented by several tower functions

# Multilayer Network

- All tower functions are similar and only differ in height and position on x-axis
- A black box takes some input and constructs a tower function
  - A network can add them up to approximate the function
- If we take the logistic function and set $w$ to a very high value, we can recover step function
  - $w$ controls the slope of the logistic function
- Can also adjust value of $b$ to control position on x-axis at which function transitions from 0 to 1
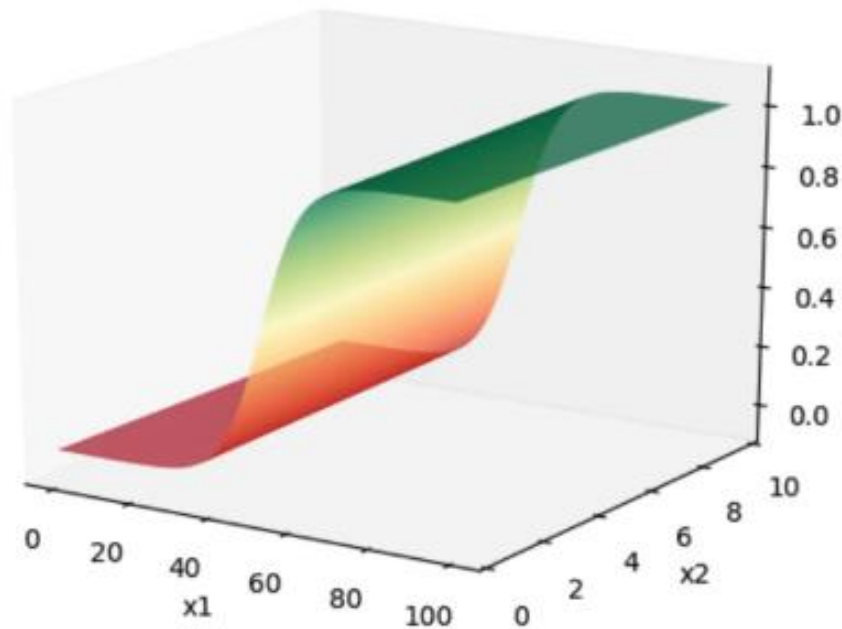
# Multilayer Network

Take two such sigmoid functions, with different $b$'s, and subtract them – will get a tower function
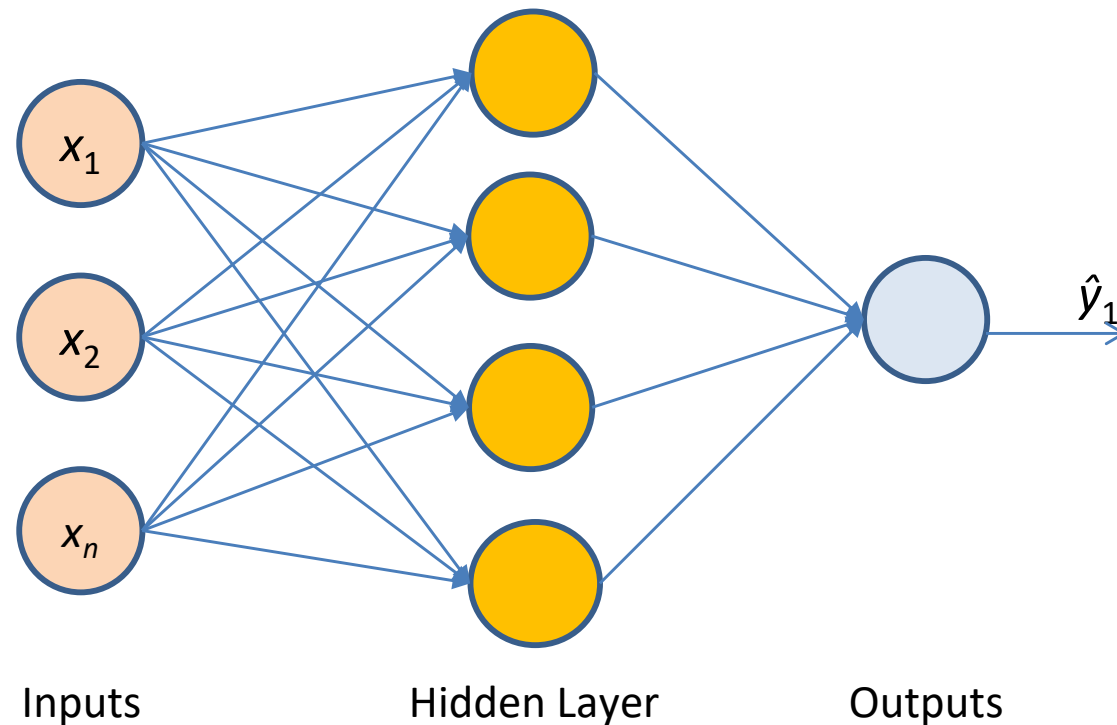
# Multilayer Network

- More input parameters??
- Ex. 2 parameters

# Single Hidden Layer Neural Network
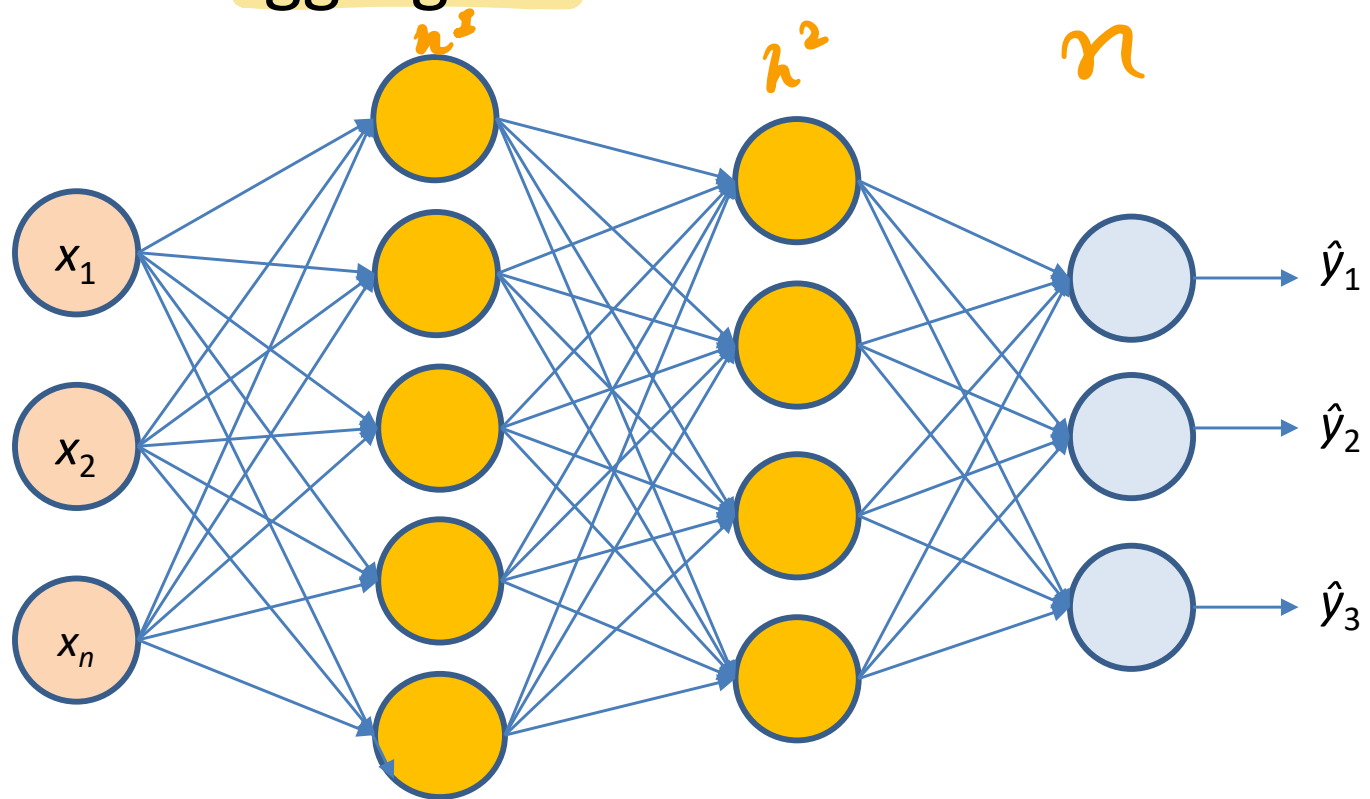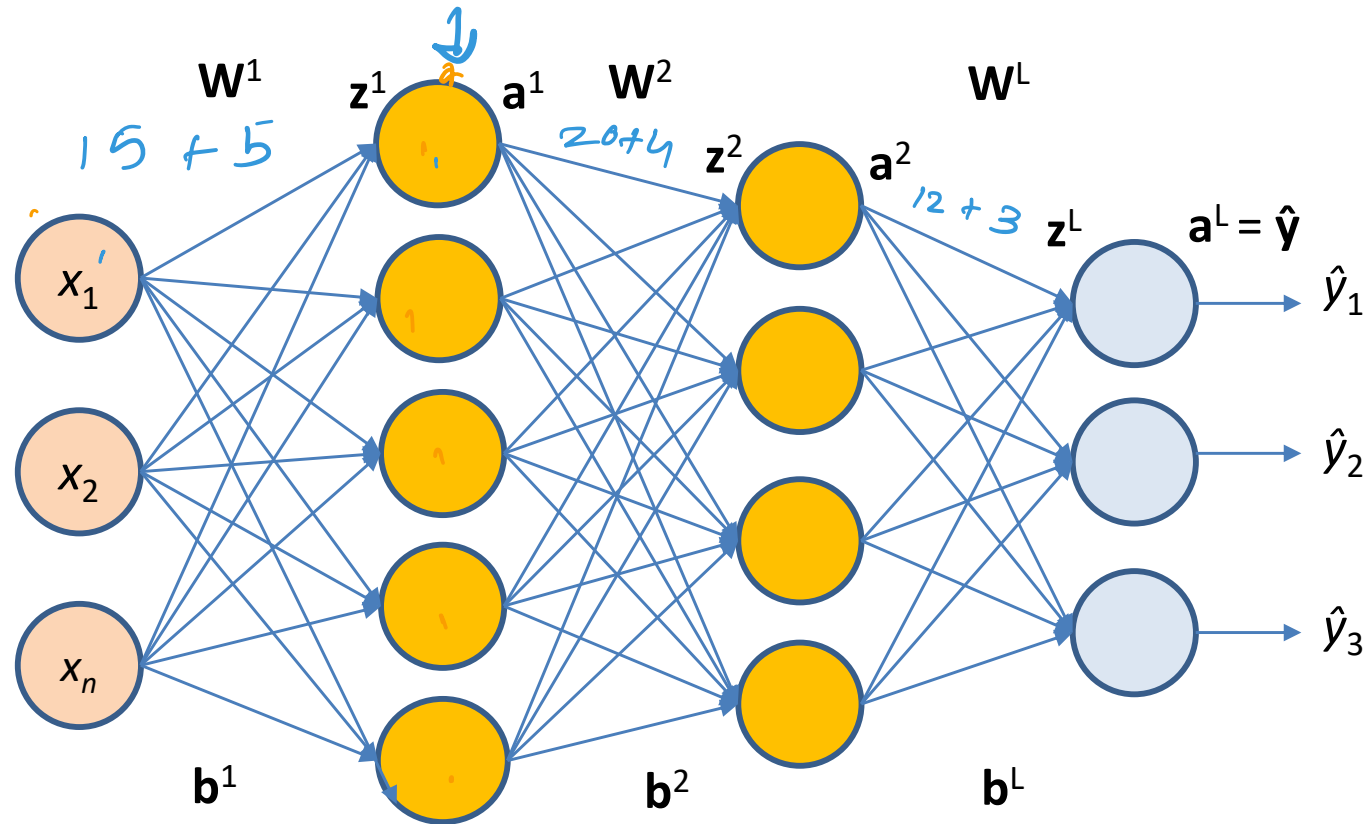


Hidden Layer: States of nodes are unobserved
Inputs are densely connected to perceptrons, hence they are called **Dense** layers or **Fully Connected** layers

# Feedforward Neural Network

- Input is an $n$-dimensional vector ($0^{th}$ layer) $\in R^n$
- Network has $L-1$ hidden layers
- 1 output layer containing $k$ neurons (ex. for $k$ classes)
- Each neuron – aggregation and activation

# Feedforward Neural Network

2 – 1

$\mathbf{W}^1$ 15 + 5 $\mathbf{z}^1$ $\mathbf{a}^1$ $\mathbf{W}^2$ 20+4 $\mathbf{z}^2$ $\mathbf{a}^2$ 12 + 3 $\mathbf{z}^L$ $\mathbf{a}^L = \hat{\mathbf{y}}$



$x_1$ $x_2$ $x_n$

$\mathbf{b}^1$ $\mathbf{b}^2$ $\mathbf{b}^L$

$\hat{y}_1$ $\hat{y}_2$ $\hat{y}_3$

Assuming $n^i$ neurons in hidden layer $h^i$, $W^i \in R^{n(i-1)*ni}$ and $b^i \in R^{ni}$ between layers $i$ -1 and $i$ for $0<i<L$

$W^L \in R^{ni*k}$ and $b^L \in R^k$ between last hidden layer and output layer

Aggregation at layer $i$ : $\mathbf{z}^i = \mathbf{W}^i \mathbf{a}^{i-1} + \mathbf{b}^i$

For first hidden layer: $\mathbf{z}^1 = \mathbf{W}^1 \mathbf{a}^0 + \mathbf{b}^1$

$$
\begin{pmatrix} z_1^{\;1} \\ z_2^{\;1} \\ z_3^{\;1} \end{pmatrix} = \begin{pmatrix} W_{11} & W_{12} & W_{13} \\ W_{21} & W_{22} & W_{23} \\ W_{31} & W_{32} & W_{33} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \\ b_3 \end{pmatrix} = \begin{pmatrix} \sum W_{1i} x_i + b_1 \\ \sum W_{2i} x_i + b_2 \\ \sum W_{3i} x_i + b_3 \end{pmatrix}
$$

Activation at layer $i$ = $g(\mathbf{z}^i)$ = $g(\mathbf{b}^i + \mathbf{W}^i \mathbf{a}^{i-1})$

For first hidden layer: $g(\mathbf{z}^1) = g(\mathbf{b}^1 + \mathbf{W}^1 \mathbf{a}^0)$

$$\begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix} = \begin{bmatrix} g(z_1) \\ g(z_2) \\ g(z_3) \end{bmatrix}$$

Eg. $g(z_1) = \sigma(z_1) = 1 / (1 + e^{-z1})$

**g: activation function (logistic, tanh, linear etc.)**

$y = \tanh$

$y' = 1 - y^2$

Handwritten annotations (orange):

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

$$\sigma'(z) = \sigma(z)(1-\sigma(z))$$

Aggregation at output layer $L = z^L = \mathbf{W^L a^{L-1} + b^L}$

$$z_1 = w_{11}x_1 + w_{12}x_2 + w_{13}x_3 + b$$

$$z_2 = w_{21}x_1 + w_{22}x_2 + w_{23}x_3 + b$$

Activation at output layer $L = \hat{\mathbf{y}} = g(z^L) = g(\mathbf{W^L a^{L-1} + b^L})$

$$\begin{bmatrix} a_1 \\ a_2 \end{bmatrix} = \begin{bmatrix} g(z_1) \\ g(z_2) \end{bmatrix}$$

# Learning parameters

In given example, dimensions of parameters:

- $\mathbf{W^1}$: $n^1 * n$ $\qquad$ $\mathbf{b^1}$: $n^1$
- $\mathbf{W^2}$: $n^2 * n^1$ $\qquad$ $\mathbf{b^2}$: $n^2$
- $\mathbf{W^L}$: $n^2 * k$ $\qquad$ $\mathbf{b^L}$: $k$

<br>

- Assuming $L$ layers and $n^i$ neurons in hidden layer $h^i$ and $k$ neurons in output layer, no. of parameters to be learned:
  - Weights: $(L-1)*(n^{i-1} * n^i) + (n*k)$ $\qquad$ for $0 < i < L$
  - Bias: $(L-1)*n^i + k$

# Learning parameters

- **Data:** $\{x_i, y_i\}$ $\qquad\qquad$ $i = 1..m$
- **Model:** $\boxed{\mathcal{w}, b}$

    $$\hat{\mathbf{y}} = f(\mathbf{x}) = g(\mathbf{W^3}g(\mathbf{W^2}g(\mathbf{W^1}\mathbf{x} + \mathbf{b^1}) + \mathbf{b^2}) + \mathbf{b^3})$$

    $$\hat{\mathbf{y}} = [\hat{y}^1 \quad \hat{y}^2 \dots \hat{y}^k]$$

- **Algorithm:** Gradient Descent with back Propagation
- **Loss/Error function:** Sum of squared error loss

    $$min \frac{1}{N}\sum_{i=1}^{m}\sum_{j=1}^{k}(\hat{y}_j^i - y_j^i) \qquad \text{for } i^{th} \text{ sample for all classes } j$$

# Learning parameters

- Gradient Descent:

    $t:=0;$

    $max\_iterations:=1000;$

    Initialize $\boldsymbol{\theta_0} := [\mathbf{W^1}_0, ... \mathbf{W^L}_0, \mathbf{b^1}_0 \ ... \ \mathbf{b^L}_0];$

    while $t{++} < max\_iterations$ do

    $\qquad \boldsymbol{\theta_{t+1}} := \boldsymbol{\theta_t} - \eta \nabla \boldsymbol{\theta_t};$

    end

where, $\nabla \theta_t = [\ \dfrac{\partial L(\theta)}{\partial W_t}, \dfrac{\partial L(\theta)}{\partial b_t}\ ]^{\mathrm{T}}$

$\nabla \theta$ composed of:

- $\nabla W^1, \nabla W^2, ... \nabla W^{L-1} \in R^{n(i-1) \times ni}\ ,\ \nabla W^L \in R^{n \times k}$
- $\nabla b^1, \nabla b^2, ... \nabla b^{L-1} \in R^{ni}\ ,\ \nabla b^L \in R^{k}$

# Loss function

- Loss function should capture how much $\hat{y}_i$ deviates from $y_i$

- $y_i \in R^n$ then squared error loss can be used:

$$L(\theta) = (1/m) * \sum (y_i - \hat{y}_i)^2$$

- Problems with squared error loss:

$$\frac{\partial L(w,b)}{\partial w} = (\hat{y} - y) * \hat{y} * (1 - \hat{y}) * x$$

  - If $y_i = 1$ and $\hat{y}_i \sim 0$, $\frac{\partial L(w,b)}{\partial w} \sim 0$     Undesirable

  - If $y_i = 0$ and $\hat{y}_i \sim 1$, $\frac{\partial L(w,b)}{\partial w} \sim 0$     Undesirable

  - Weight updation becomes very slow

# Loss function

- Information content (IC):
  - Events with high probability have low information content
    - "The sun will rise tomorrow"
  - Events with low probability have high information content
    - "There will be a cyclone tomorrow"
- $IC(A) = -\log_2(p(A))$
- Entropy: Expected information content $= \sum p_i * IC(i)$

$$= - \sum p_i \log_2(p_i)$$

# Loss function

Entropy: $y_i$ = [0    1    0    0]        //Team B wins game

$\hat{y}_i$ = [0.2   0.1   0.4   0.3]      //Our prediction

10K   5K   8K   1K     //Profit for each team win

Expected profit??

- Entropy: Expected information content = $\sum p_i \, IC(i)$

$= -\sum p_i \log_2(p_i)$

Actual    predicted.

# Loss function

- Cross-entropy: gives a measure on how close a predicted distribution is to a true distribution
  - True distribution $p_i$, Estimated distribution $q_i$
  - Estimated information content = $-\sum p_i \log_2(q_i)$
  - Capture difference between two probability distributions
  - If prediction is close to actual, cross entropy will be low

$$L(\theta) = -\sum y_c \log_2(\hat{y}_c) \qquad \text{for all } k \text{ classes}$$

$$y_c = 1 \qquad \text{if } c = t \text{ (true class)}$$

$$= 0 \qquad \text{otherwise}$$

$$L(\theta) = -\log_2(\hat{y}_t)$$

# Loss function

- Objective function for classification:
  - Cross-entropy Loss

    minimize: $L(\theta) = -\log_2(\hat{y}_t)$

$$\frac{\partial L}{\partial w} = -\frac{1}{\hat{y}_t} \frac{\partial \hat{y}}{\partial w} \leftarrow$$

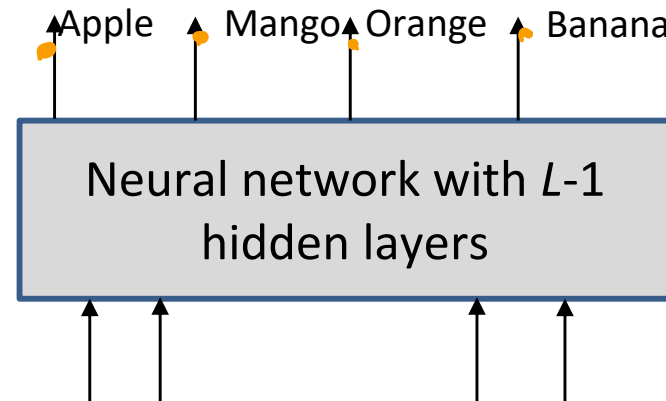$\hat{y}_t$: predicted probability of correct event

$\log_2(\hat{y}_t)$: probability that $x$ belongs to $t^{th}$ class, log-likelihood of data

# Output Activation Function

- Output activation function:
  - Sum of outputs should be 1
  - $\hat{y}$ should be a probability distribution
  - Sigmoid – probabilities will be $0<p<1$ but sum not equal to 1

$\sigma \longrightarrow$

Softman

$y_i = \{1 \qquad 0 \qquad 0 \qquad 0\}$

Apple    Mango  Orange    Banana

Neural network with $L$-1 hidden layers

$$\dfrac{e^{\top z_1}}{\displaystyle\sum_{i=1}^{k} e^{\top z_i}}$$

**Classification problem**

**apple**

# Output Activation Function

$$\frac{e^z}{z}$$

- Softmax function

  $$z^L = b^L + W^L \, a^{L-1}$$

  $$\hat{y} = g(z^L_j) = e^z{}_j / \sum e^z{}_j \qquad\qquad \text{for } j = 1..k$$

  $z^L_j$ is $j^{\text{th}}$ element of $z^L$

- Example: $z^L = [10 \quad 20 \quad -30]$

  $\hat{y} = [e^{10}/(e^{10} + e^{20} + e^{-30}) \qquad e^{20}/(e^{10} + e^{20} + e^{-30}) \qquad e^{-30}/(e^{10} + e^{20} + e^{-30}) \ ]$

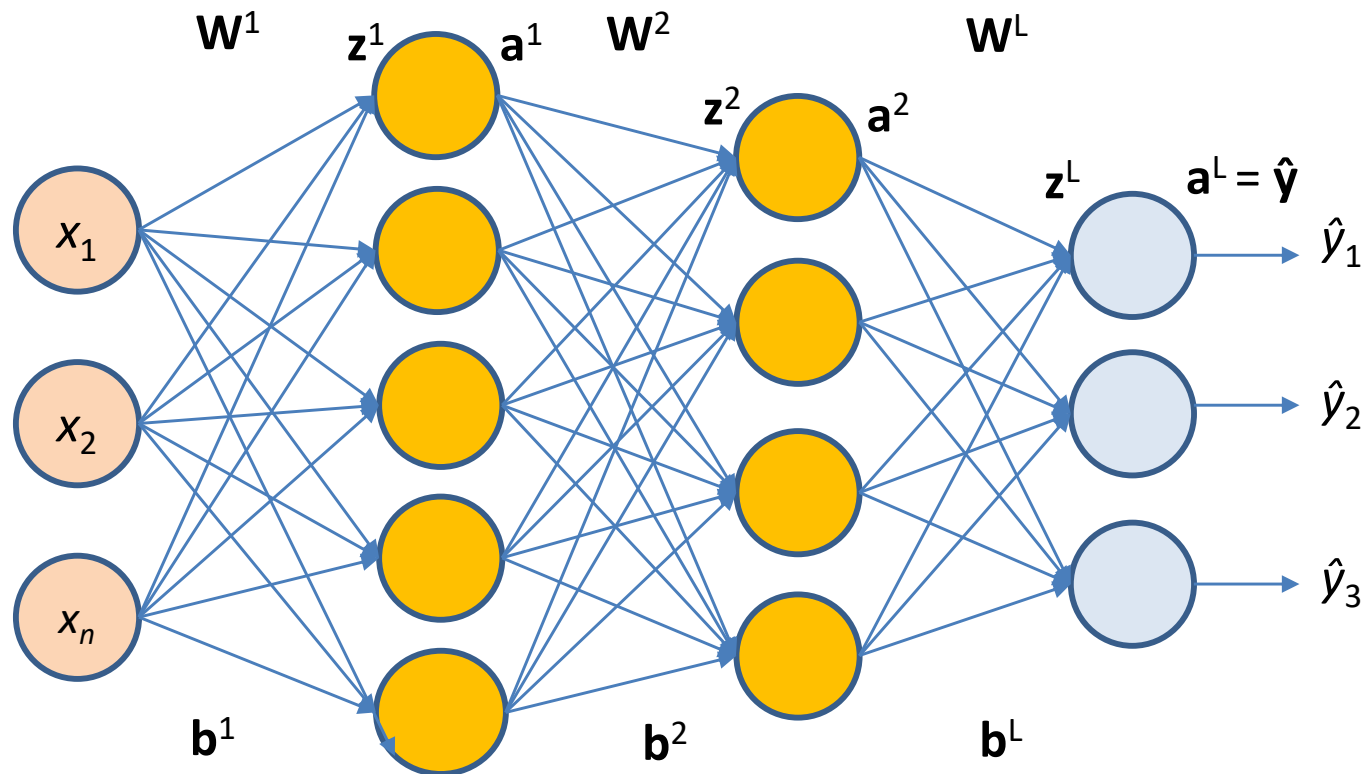NOTE: Exponent converts –ve values to +ve values

# Loss function

| | Outputs | |
|---|---|---|
| | **Real values** | **Probabilities** |
| Output activation | Linear | Softmax |
| Loss function | Squared error | Cross-entropy |

# Backpropagation

How to compute $\nabla\theta$ composed of:

$\nabla W^1, \nabla W^2, \dots \nabla W^{L-1} \in R^{n \times n}$ , $\nabla W^L \in R^{n \times k}$

$\nabla b^1, \nabla b^2, \dots \nabla b^{L-1} \in R^n$ , $\nabla b^L \in R^k$
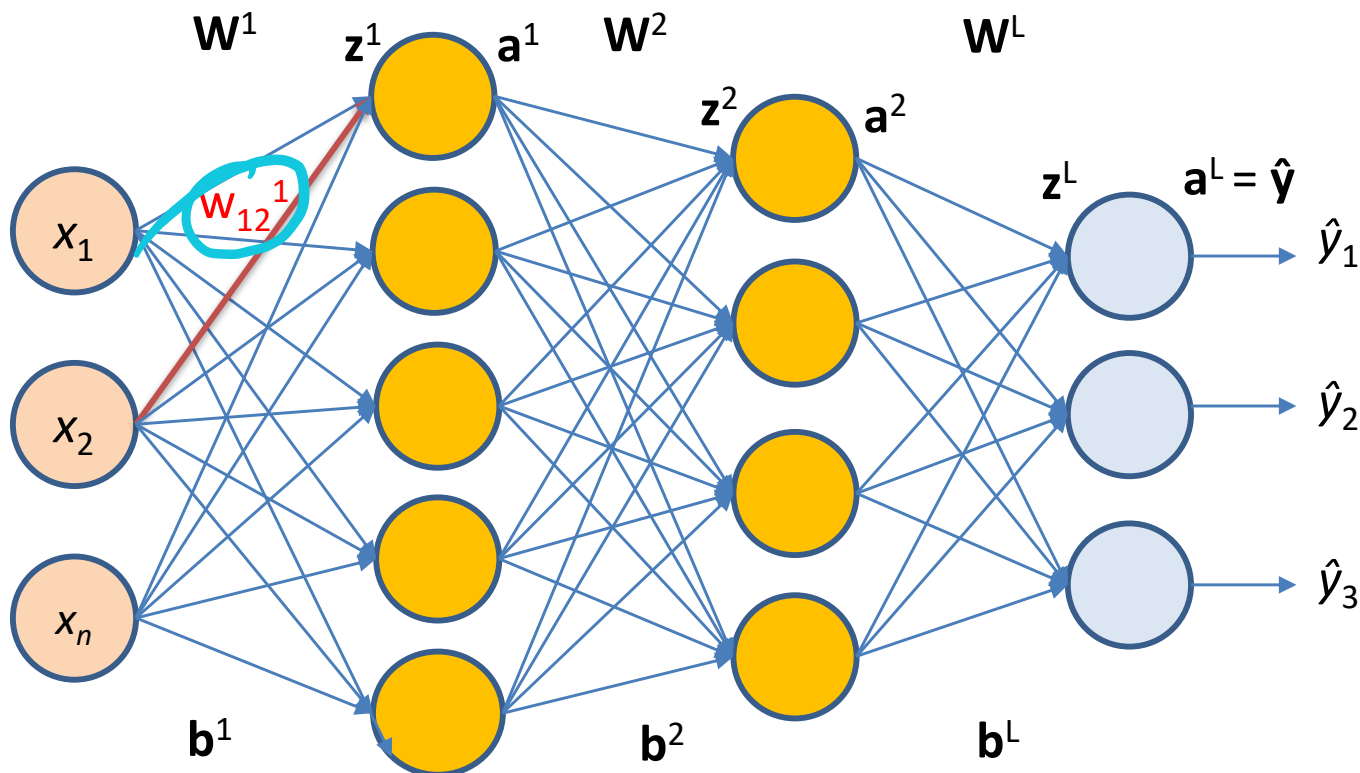
$$\frac{\partial L}{\partial w_{12}} = \frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial z_L} \times \frac{\partial z_L}{\partial a^2} \times \frac{\partial a^2}{\partial z^2} \times \frac{\partial z^2}{\partial a^1} \times \frac{\partial a^1}{\partial z^1} \times \frac{\partial z_1}{\partial w_{12}}$$
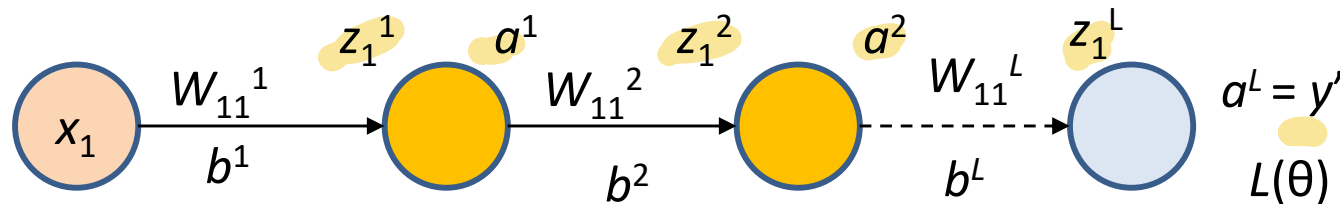
# Backpropagation

Assuming classification problem, $L(\theta) = -\log_2(\hat{y}_t)$

- To learn weight $w_{12}^1$ use SGD and compute $\frac{\partial L(w,b)}{\partial W_{12}}$

# Backpropagation

Assume a deep thin network, who is responsible for the loss??



Find derivative by chain rule:

$$\frac{\partial L(\theta)}{\partial W_{11}^1} = \frac{\partial L(\theta)}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial z_1^L} * \frac{\partial z_1^L}{\partial a_1^2} * \frac{\partial a_1^2}{\partial z_1^2} * \frac{\partial z_1^2}{\partial a_1^1} * \frac{\partial a_1^1}{\partial z_1^1} * \frac{\partial z_1^1}{\partial W_{11}^1}$$

Output layer      Previous hidden layer      Previous hidden layer      Weights
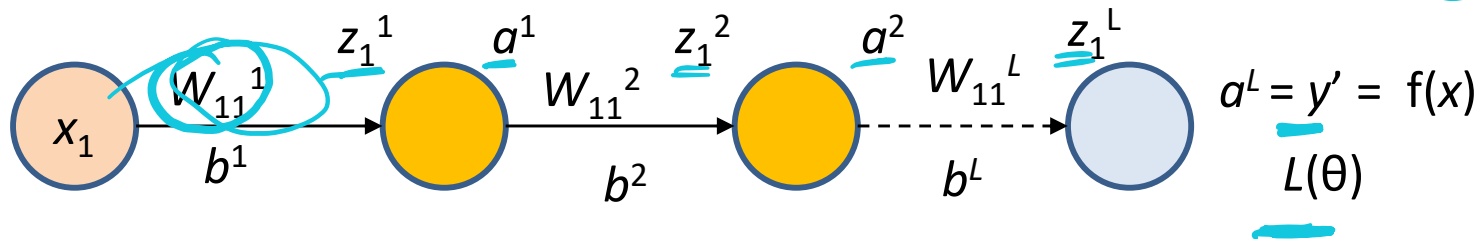
**If we change $W_{11}$, how much does the loss change**

# Backpropagation

Assume a deep thin network

SSE

$\text{Cross Entropy} = -\sum_k y_k \log_2(\hat{y}_i)$



Find derivative by chain rule:

$$\frac{\partial L(\theta)}{\partial W_{11}^1} = \frac{\partial L(\theta)}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial z_1^L} * \frac{\partial z_1^L}{\partial a_1^2} * \frac{\partial a_1^2}{\partial z_1^2} * \frac{\partial z_1^2}{\partial a_1^1} * \frac{\partial a_1^1}{\partial z_1^1} * \frac{\partial z_1^1}{\partial W_{11}^1}$$

$$\frac{\partial L(\theta)}{\partial W_{11}^2} = \frac{\partial L(\theta)}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial z_1^L} * \frac{\partial z_1^L}{\partial a_1^2} * \frac{\partial a_1^2}{\partial z_1^2} * \frac{\partial z_1^2}{\partial W_{11}^2}$$

$$\frac{\partial L(\theta)}{\partial W_{11}^L} = \frac{\partial L(\theta)}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial z_1^L} * \frac{\partial z_1^L}{\partial W_{11}^L}$$

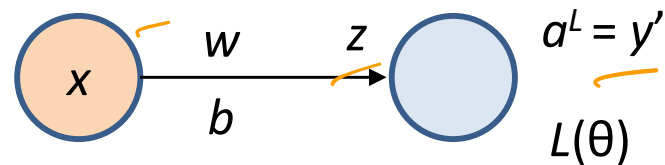$$L(w,b) = -y \log \hat{y} - (1-y) \log(1-\hat{y})$$

$$\frac{\partial L}{\partial \hat{w}} = \left[ \frac{-y}{\hat{y}} + \frac{(1-y)}{(1-\hat{y})} \right] \hat{y}(1-\hat{y}) x$$

$$\frac{\partial y}{\partial z} \times \partial$$

$$= (\hat{y} - y) \_ \_ \_$$

# BACKPROPAGATION WITH SIGMOID OUTPUT ACTIVATION & BINARY CROSS-ENTROPY LOSS

# Backpropagation

$x$ $\xrightarrow[b]{w}$ $z$ $\to$ $a^L = y'$

$L(\theta)$

Assuming binary cross-entropy function

$L = -y \log \hat{y} - (1 - y) \log (1 - \hat{y})$

$$\frac{\partial L(\theta)}{\partial \hat{y}} = \frac{-y}{\hat{y}} + \frac{1 - y}{1 - \hat{y}}$$

$$\frac{\partial L(\theta)}{\partial z} = \frac{\partial L(\theta)}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial z} = \left( \frac{-y}{\hat{y}} + \frac{1-y}{1-\hat{y}} \right) * \hat{y}(1 - \hat{y}) = \hat{y} - y$$

$$\frac{\partial L(\theta)}{\partial w} = \frac{\partial L(\theta)}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial z} * \frac{\partial \hat{z}}{\partial w} = (\hat{y} - y) * x$$
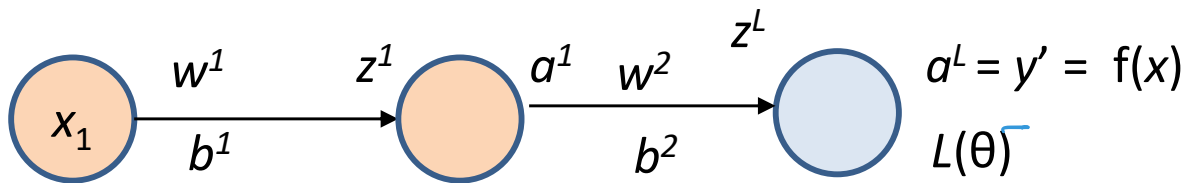
Loss $\to$ Cross entropy

Activation $\to$ Sigmoid

$$\frac{\partial L(\theta)}{\partial \hat{y}} = \frac{-y}{\hat{y}} + \frac{1-y}{1-\hat{y}}$$

$$\frac{\partial L(\theta)}{\partial z^2} = \frac{\partial L(\theta)}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial z^2} = \left(\frac{-y}{\hat{y}} + \frac{1-y}{1-\hat{y}}\right) * \hat{y}(1-\hat{y}) = \hat{y} - y \qquad :\delta^L = \frac{\partial L}{\partial \hat{y}} * \sigma'(z^L)$$

$$\frac{\partial L(\theta)}{\partial w^2} = \frac{\partial L(\theta)}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial z^2} * \frac{\partial z^2}{\partial w^2} = (\hat{y} - y) * a^1 \qquad :\delta^L * a^1$$

$$\frac{\partial L(\theta)}{\partial a^1} = \frac{\partial L(\theta)}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial z^2} * \frac{\partial z^2}{\partial a^1} = (\hat{y} - y) * w^2 \qquad : \delta^L * w^2$$

$$\frac{\partial L(\theta)}{\partial z^1} = \frac{\partial L(\theta)}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial z^2} * \frac{\partial z^2}{\partial a^1} * \frac{\partial a^1}{\partial z^1} = (\hat{y} - y) * w^2 * a^1(1-a^1) \qquad :\delta^1 = \delta^L w^2 * \sigma'(z^1)$$

$$\frac{\partial L(\theta)}{\partial w^1} = \frac{\partial L(\theta)}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial z^2} * \frac{\partial z^2}{\partial a^1} * \frac{\partial a^1}{\partial z^1} * \frac{\partial z^1}{\partial w^1} = (\hat{y} - y) * w^2 * a^1(1-a^1) * x \qquad :\delta^1 * x$$

Handwritten annotations:

$$z^2 = w^2 a_1 + b \qquad a^1 = \frac{1}{1+e^{-z^1}}$$

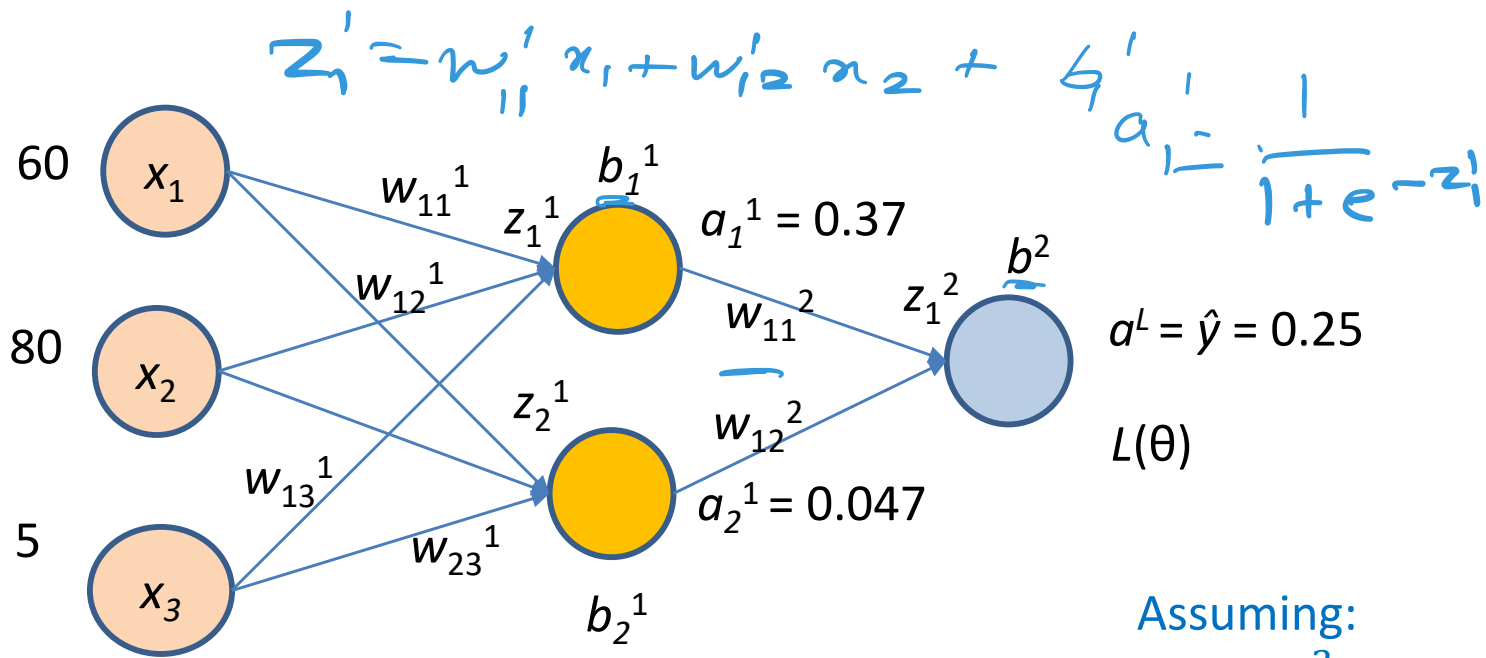$$\hat{y} = \frac{1}{1+e^{-z_2}} \qquad z^1 = w_1 a_1 + b_1$$

$$w2$$

$$z^2 = w^2 a_1 + b^2$$

# Backpropagation Equations

- $\delta^L = \dfrac{\partial L}{\partial \hat{y}} \odot \sigma'(z^L)$

- $\delta^l = \left(\left(w^{l+1}\right)^T \delta^{l+1}\right) \odot \sigma'\left(z^l\right)$

- $\dfrac{\partial L}{\partial b_j^l} = \delta_j^l$

- $\dfrac{\partial L}{\partial w_{jk}^l} = a_k^{l-1}\delta_j^l$

$$z_1' = w_{11}' x_1 + w_{12}' x_2 + b_1'$$

$$a_1' = \frac{1}{1+e^{-z_1'}}$$



60  $x_1$

$w_{11}^1$  $z_1^1$  $b_1^1$  $a_1^1 = 0.37$

$w_{12}^1$

80  $x_2$  $w_{11}^2$  $z_1^2$  $b^2$

$z_2^1$  $a^L = \hat{y} = 0.25$

$w_{13}^1$  $w_{12}^2$

5  $a_2^1 = 0.047$  $L(\theta)$

$x_3$  $w_{23}^1$

$b_2^1$

$y=1$

Assuming:
Initial $w_{11}^2 = 12$, $w_{11}^1 = 0.1$
$\eta = 0.01$

$$\frac{\partial L}{\partial w_{11}^2} = \frac{\partial L}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial z_1^2} * \frac{\partial z_1^2}{\partial w_{11}^2}$$

$$\frac{\partial L}{\partial \hat{y}} \times \frac{\partial \hat{y}}{\partial z_1^2} \times \frac{\partial z_1^2}{\partial w_{11}^2}$$

During forward propagation:
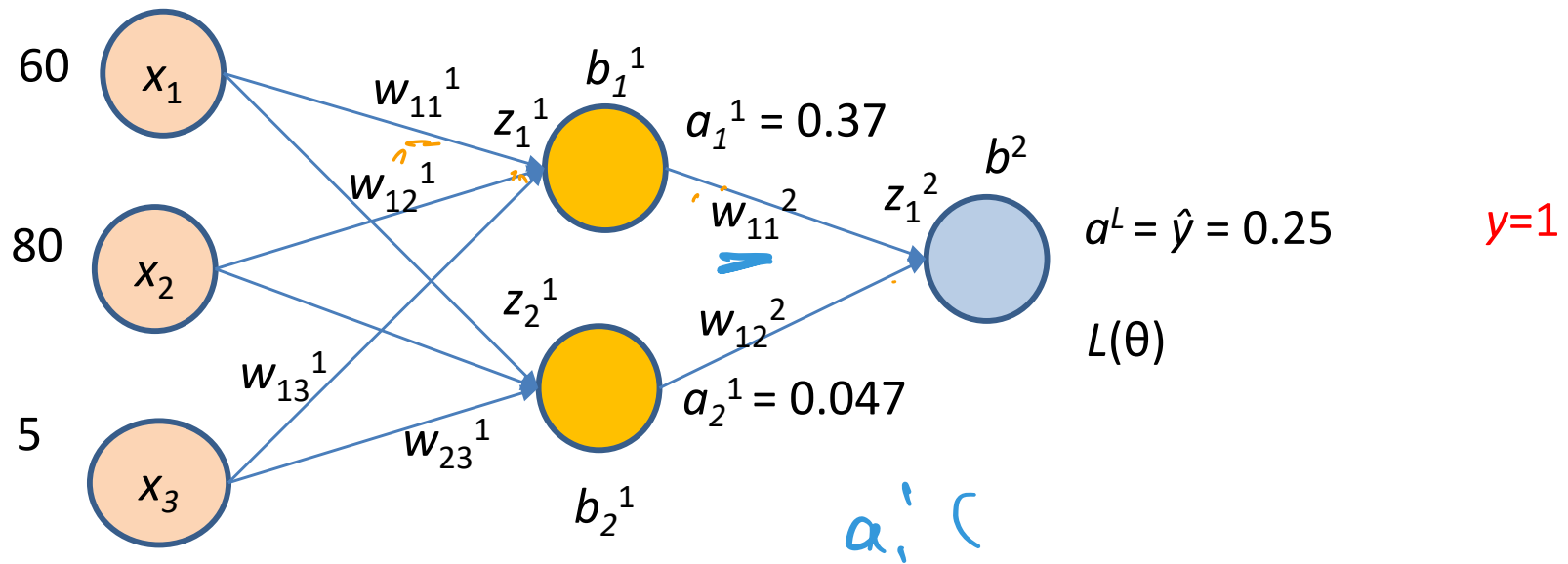$a_1^1 = 0.37$, $z_1^1 = 0.5$

$$\frac{\partial L(\theta)}{\partial \hat{y}} = \frac{-y}{\hat{y}} + \frac{1-y}{1-\hat{y}} = \frac{-1}{0.25} = -4$$

$$\frac{\partial L(\theta)}{\partial z_1^2} = \frac{\partial L(\theta)}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial z_1^2} = \hat{y} - y = 0.25 - 1 = -0.75$$

$$\left( \frac{-y}{\hat{y}} + \frac{1-y}{1-\hat{y}} \right) \hat{y}(1-\hat{y}) \times a_1^1$$

$$\frac{\partial L(\theta)}{\partial w_{11}^2} = \frac{\partial L(\theta)}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial z_1^2} * \frac{\partial z}{\partial w_{11}^2} = -0.75 * a_1^1 = -0.75 * 0.37 = -0.2775$$

$$w_{11}^{2*} = w_{11}^2 - \eta * \frac{\partial L}{\partial w_{11}^2} = 12 - 0.01 * (-0.2775) = 12.0028$$

$$\frac{\partial L}{\partial w_{11}^1} = \frac{\partial L}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial z_1^2} * \frac{\partial z_1^2}{\partial a_1^1} * \frac{\partial a_1^1}{\partial z_1^1} * \frac{\partial z_1^1}{\partial w_{11}^1}$$

$$= -0.75 * w_{11}^2 * a_1^1(1 - a_1^1) * x_1$$

= -0.75 * 12 * 0.37 (1- 0.37 )*60

= -125.874

$$w_{11}^{1*} = w_{11}^1 - \eta * \frac{\partial L}{\partial w_{11}^1} = 0.1 - 0.01 * (-125.874) = 1.35$$

**Assignment: Compute** $\frac{\partial L}{\partial w_{13}^1}$

# BACKPROPAGATION WITH SOFTMAX OUTPUT ACTIVATION & CROSS-ENTROPY LOSS

# Backpropagation

- Computing gradients
  - Gradients w.r.t. output units
  - Gradients w.r.t. hidden units
  - Gradients w.r.t. weights and biases

$$\frac{\partial L(\theta)}{\partial W_{11}^1} = \frac{\partial L(\theta)}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial z_1^L} * \frac{\partial z_1^L}{\partial a_1^2} * \frac{\partial a_1^2}{\partial z_1^2} * \frac{\partial z_1^2}{\partial a_1^1} * \frac{\partial a_1^1}{\partial z_1^1} * \frac{\partial z_1^1}{\partial W_{11}^1}$$

# Gradients w.r.t. output units

- Assuming softmax activation and cross entropy loss at output layer for *k* classes:

$$L(\theta) = -\log_2(\hat{y}_t) \quad t\text{: true class label}$$

$$\frac{\partial L(\theta)}{\partial \hat{y}_i} = \frac{\partial(-log\hat{y}_t)}{\partial \hat{y}_i} \qquad \text{for } i = 1..k$$

$$= -\frac{1}{\hat{y}_t} \qquad \text{if } i = t$$

$$= 0 \qquad \text{otherwise}$$

# Gradients w.r.t. output units

$$\frac{\partial L(\theta)}{\partial \hat{y}_i} = -\frac{1}{\hat{y}_t} \qquad \text{if } i = t$$
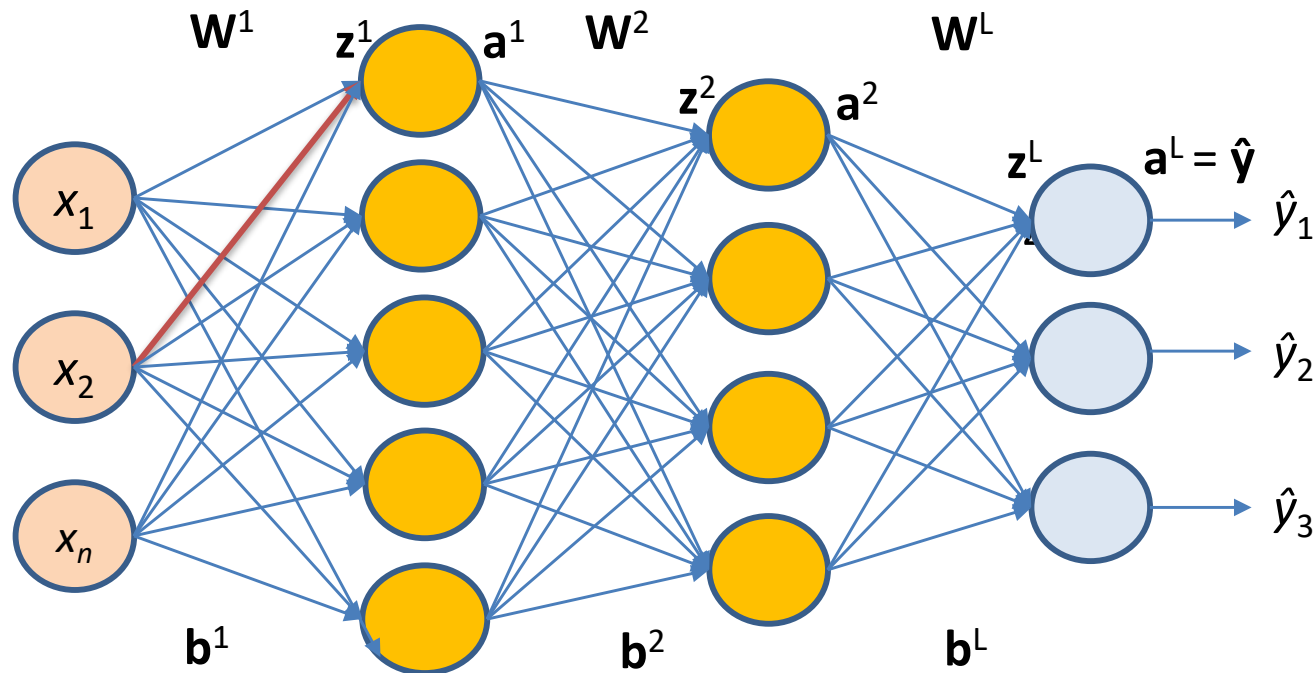
$$= \quad 0 \qquad \qquad \text{otherwise}$$

$$\frac{\partial L(\theta)}{\partial \hat{y}} = \begin{pmatrix} \dfrac{\partial L(\theta)}{\partial \hat{y}_1} \\[2ex] \dfrac{\partial L(\theta)}{\partial \hat{y}_2} \\[2ex] .... \\[2ex] \dfrac{\partial L(\theta)}{\partial \hat{y}_k} \end{pmatrix} = -\frac{1}{\hat{y}_t} * \begin{pmatrix} 1 \text{ if } t = 1 \\[2ex] 1 \text{ if } t = 2 \\[2ex] ... \\[2ex] 1 \text{ if } t = k \end{pmatrix}$$

$$= -\frac{1}{\hat{y}_t} * e(t)$$

$e(t)$: One hot $k$-dimensional vector whose $t^{\text{th}}$ entry is 1, others are 0

# Gradients w.r.t. output units

$$\frac{\partial L(\theta)}{\partial z_i^L} = \frac{\partial(-\log \hat{y}_t)}{\partial z_i^L} = -\frac{\partial(-\log \hat{y}_t)}{\partial \hat{y}_t} * \frac{\partial \hat{y}_t}{\partial z_i^L} = -\frac{1}{\hat{y}_t} * \frac{\partial \hat{y}_t}{\partial z_i^L}$$

$$\hat{y}_t = \frac{\exp(z_t^L)}{\sum \exp(z_i^L)} \longrightarrow \quad \hat{y}_t \text{ depends on } z_i^L$$

# Gradients w.r.t. output units

$$\frac{\partial(-\log \hat{y}_t)}{\partial z_i^L} = -\frac{1}{\hat{y}_t} * \frac{\partial \hat{y}_t}{\partial z_i^L}$$

$z^L = [z_1{}^L \quad z_2{}^L \quad \dots \quad z_k{}^L]$

$\hat{y} = \text{softmax}(z^L) = [\hat{y}_1{}^L \quad \hat{y}_2{}^L \quad \dots \hat{y}_t{}^L \quad \dots \hat{y}_k{}^L]$

$\hat{y}_t$: $t^{\text{th}}$ entry of $\hat{y} = \text{softmax}(z_t{}^L) = \dfrac{\exp(z_t^L)}{\sum \exp(z_i^L)}$

$$\frac{\partial(-\log \hat{y}_t)}{\partial z_i^L} = -\frac{1}{\hat{y}_t} * \frac{\partial \hat{y}_t}{\partial z_i^L}$$

$$= -\frac{1}{\hat{y}_t} * \frac{\partial(softmax(z_t^L)}{\partial z_i^L}$$

# Gradients w.r.t. output units

$$\frac{\partial(-\log \hat{y}_t)}{\partial z_i^L} = -\frac{1}{\hat{y}_t} * \frac{\partial(softmax(z_t^L)}{\partial z_i^L} = -\frac{1}{\hat{y}_t} * \frac{\partial}{\partial z_i^L}\left(\frac{\exp(z_t^L)}{\sum \exp(z_i^L)}\right)$$

$$\frac{\partial(-\log \hat{y}_t)}{\partial z_i^L} = -\frac{1}{\hat{y}_t} * \frac{\frac{\partial\left(\exp(z_t^L)\right)}{\partial z_i^L}}{\sum \exp(z_i^L)} - \frac{\exp(z_t^L)*\frac{\partial\left(\sum \exp(z_i^L)\right)}{\partial z_i^L}}{(\sum \exp(z_i^L))^2}$$

$$= -\frac{1}{\hat{y}_t} * \left(\frac{1_{(t=i)}\exp(z_t^L)}{\sum \exp(z_i^L)} - \frac{\exp(z_t^L)*\exp(z_i^L)}{(\sum \exp(z_i^L))^2}\right)$$

$$= -\frac{1}{\hat{y}_t} * \left(1_{(t=i)}\text{softmax}(z_t^L) - \text{softmax}(z_t^L) * \text{softmax}(z_i^L)\right)$$

$$= -\frac{1}{\hat{y}_t} * \left(1_{(t=i)}\hat{y}_t - \hat{y}_t \hat{y}_i\right) = -\left(1_{(t=i)} - \hat{y}_i\right)$$

$$\frac{\partial}{\partial x}\frac{f(x)}{g(x)} = \frac{g(x)*\frac{\partial(f(x))}{\partial x} - f(x)*\frac{\partial(g(x))}{\partial x}}{(g(x))^2} = \frac{\frac{\partial(f(x))}{\partial x}}{g(x)} - \frac{f(x)*\frac{\partial(g(x))}{\partial x}}{(g(x))^2}$$

# Gradients w.r.t. output units

$$\frac{\partial L(\theta)}{\partial z_i^L} = -(1_{(t=i)} - \hat{y}_i)$$

Gradient w.r.t. vector **z**$^L$:
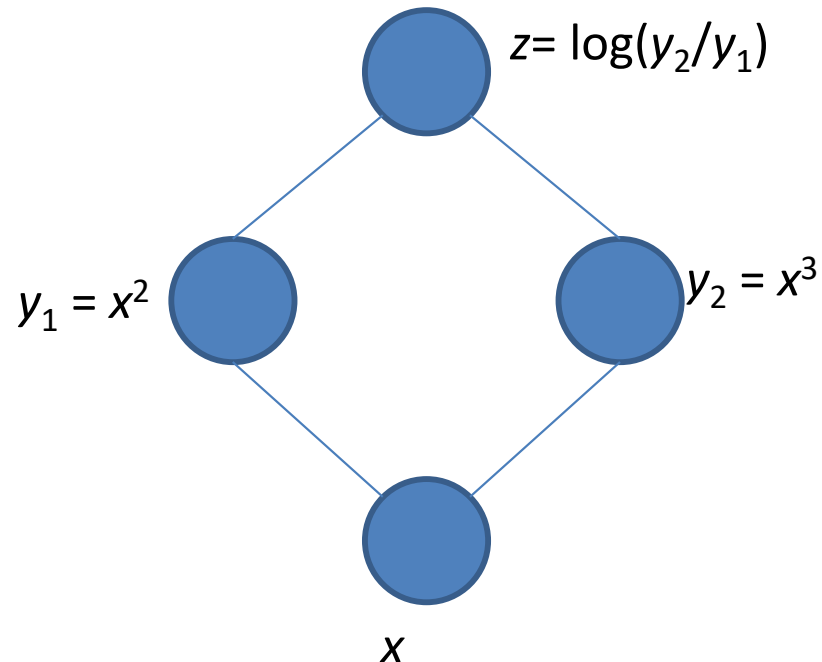
$$\frac{\partial L(\theta)}{\partial z^L} = \begin{pmatrix} \dfrac{\partial L(\theta)}{\partial z_1^L} \\ \dfrac{\partial L(\theta)}{\partial z_2^L} \\ \dots \\ \dfrac{\partial L(\theta)}{\partial z_k^L} \end{pmatrix} = \begin{pmatrix} -(1_{(t=1)} - \hat{y}_1) \\ -(1_{(t=2)} - \hat{y}_2) \\ \dots \\ -(1_{(t=k)} - \hat{y}_k) \end{pmatrix} = -(e(t) - \hat{y})$$

<span style="color:red">$e(t)$: One hot $k$-dimensional vector whose $t$th entry is 1, others are 0</span>

# Gradients w.r.t. hidden units

**Example**
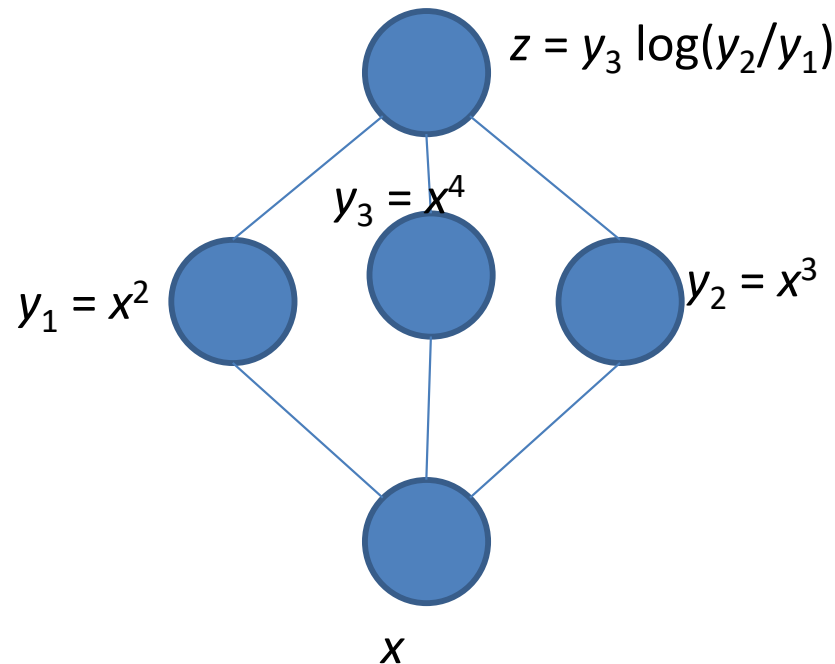
$$\frac{\partial z}{\partial x} = \frac{\partial z}{\partial y_1} * \frac{\partial y_1}{\partial x} + \frac{\partial z}{\partial y_2} * \frac{\partial y_2}{\partial x}$$

$$= \sum \frac{\partial z}{\partial y_i} * \frac{\partial y_i}{\partial x} \qquad \text{for } i = 1,2$$



$z = \log(y_2/y_1)$
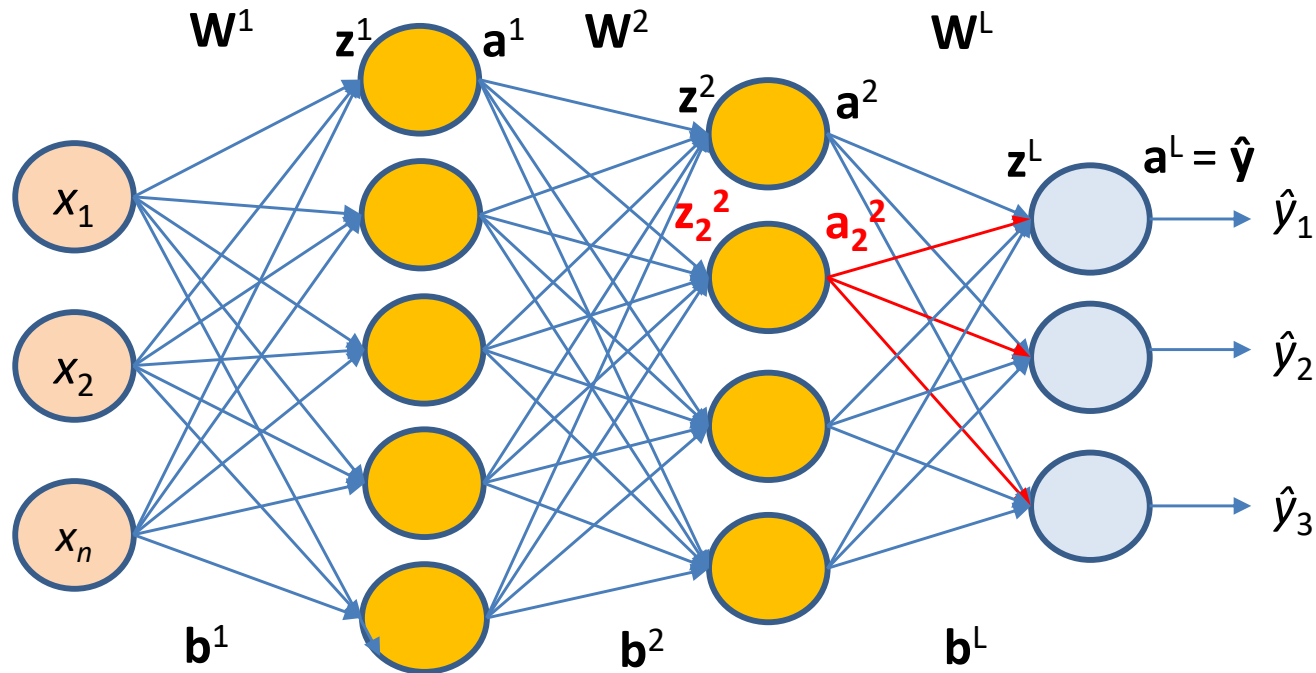
$y_1 = x^2$

$y_2 = x^3$

$x$

# Gradients w.r.t. hidden units

**Example**

$$\frac{\partial z}{\partial x} = \sum \frac{\partial z}{\partial y_i} * \frac{\partial y_i}{\partial x}$$

for $i = 1,2,3$



$z = y_3 \log(y_2/y_1)$

$y_3 = x^4$

$y_1 = x^2$

$y_2 = x^3$

$x$

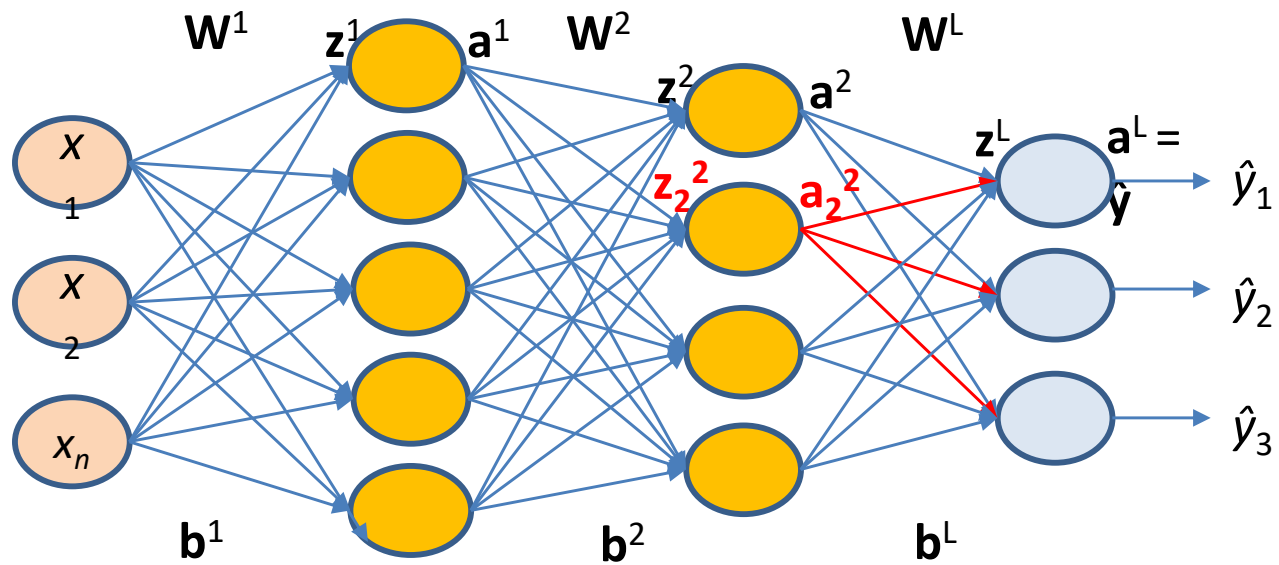**Chain rule across multiple paths**

# Gradients w.r.t. hidden units



$$z^{l+1} = W^{l+1} a_j^l + b^{l+1}$$

$$\frac{\partial L(\theta)}{\partial a_j^l} = \sum_{p=1}^k \frac{\partial L(\theta)}{\partial z_p^{l+1}} * \frac{\partial z_p^{l+1}}{\partial a_j^l}$$

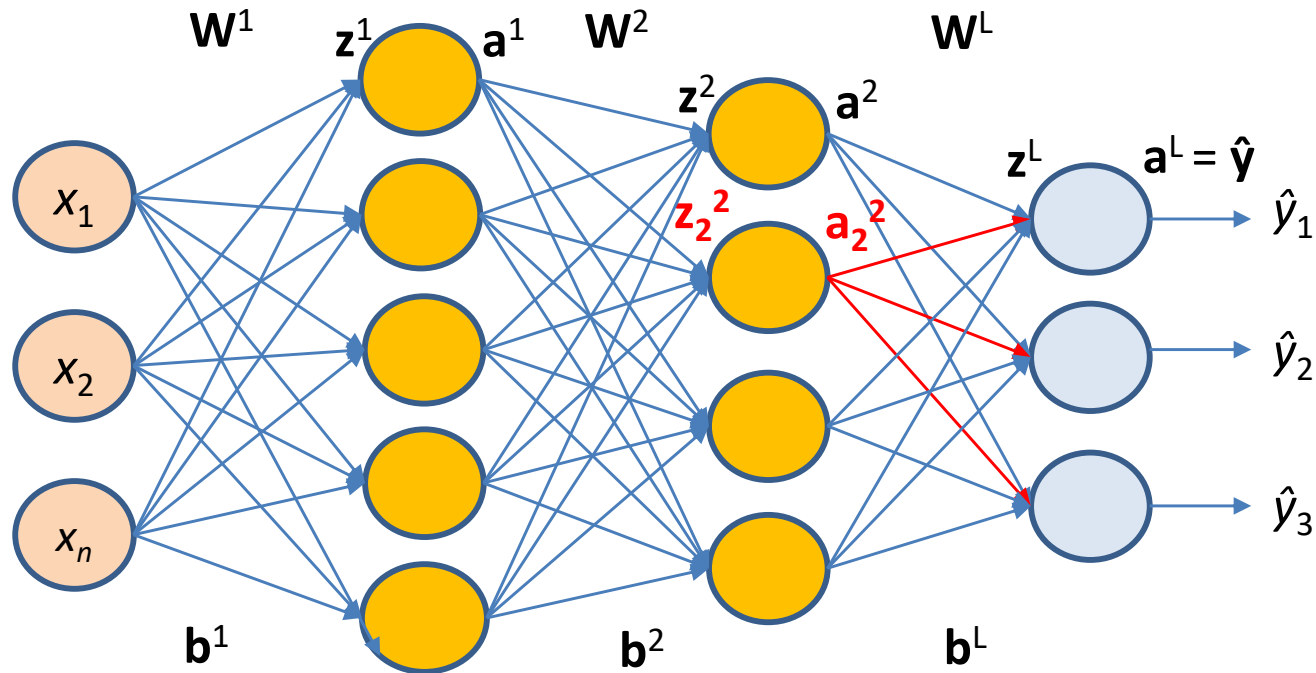**k** possible paths from output to hidden unit

$$z^{l+1} = W^{l+1} a^l_j + b^{l+1}$$

$$\begin{bmatrix} z^L_1 \\ z^L_2 \\ z^L_3 \end{bmatrix} = \begin{bmatrix} w^L_{11} & w^L_{12} & w^L_{13} & w^L_{14} \\ w^L_{21} & w^L_{22} & w^L_{23} & w^L_{24} \\ w^L_{31} & w^L_{32} & w^L_{33} & w^L_{34} \end{bmatrix} \begin{bmatrix} a^2_1 \\ a^2_2 \\ a^2_3 \\ a^2_4 \end{bmatrix} + \begin{bmatrix} b^L_1 \\ b^L_2 \\ b^L_3 \end{bmatrix}$$

$$z^L_1 = W^L_{11} a^2_1 + W^L_{12} a^2_2 + W^L_{13} a^2_3 + W^L_{14} a^2_4 + b^L_1$$

$$\frac{\partial z^L_1}{\partial a^2_2} = w^L_{12}$$

# Gradients w.r.t. hidden units



$$z^{l+1} = W^{l+1}a_j^l + b^{l+1}$$

$$\frac{\partial L(\theta)}{\partial a_j^l} = \sum_{p=1}^{k} \frac{\partial L(\theta)}{\partial z_p^{l+1}} * \frac{\partial z_p^{l+1}}{\partial a_j^l}$$

$$\frac{\partial L(\theta)}{\partial a_j^l} = \sum_{p=1}^{k} \frac{\partial L(\theta)}{\partial z_p^{l+1}} * W_{p,j}^{l+1}$$

# Gradients w.r.t. hidden units

$$\frac{\partial L(\theta)}{\partial a_j^l} = \sum_{p=1}^{k} \frac{\partial L(\theta)}{\partial z_p^{l+1}} * \frac{\partial z_p^{l+1}}{\partial a_j^l}$$

$$\frac{\partial L(\theta)}{\partial a_j^l} = \sum_{p=1}^{k} \frac{\partial L(\theta)}{\partial z_p^{l+1}} * W_{p,j}^{l+1}$$

① Vision

② Unsuper

③ Adv. Analytic

$$\frac{\partial L(\theta)}{\partial z^{l+1}} = \begin{pmatrix} \dfrac{\partial L(\theta)}{\partial z_1^{l+1}} \\ \dfrac{\partial L(\theta)}{\partial z_2^{l+1}} \\ \dots \\ \dfrac{\partial L(\theta)}{\partial z_k^{l+1}} \end{pmatrix} \qquad W_{.,j}^{l+1} = \begin{pmatrix} W_{1,j}^{l+1} \\ W_{2,j}^{l+1} \\ \dots \\ W_{k,j}^{l+1} \end{pmatrix}$$

$W_{.,j}^{l+1}$ is the $j^{\text{th}}$ column of $W^{l+1}$

$$(W_{.,j}^{l+1})^{\top} \frac{\partial L(\theta)}{\partial z^{l+1}} = \sum_{p=1}^{k} \frac{\partial L(\theta)}{\partial z_p^{l+1}} * W_{p,j}^{l+1}$$

# Gradients w.r.t. hidden units

$$\frac{\partial L(\theta)}{\partial a_j^l} = (W_{\cdot,j}^{l+1})^\mathsf{T} \frac{\partial L(\theta)}{\partial z^{l+1}}$$

$$\frac{\partial L(\theta)}{\partial a_j^l} = \begin{pmatrix} \frac{\partial L(\theta)}{\partial a_1^l} \\ \frac{\partial L(\theta)}{\partial a_2^l} \\ \dots \\ \frac{\partial L(\theta)}{\partial a_n^l} \end{pmatrix} = \begin{pmatrix} (W_{\cdot,1}^{l+1})^\mathsf{T} \frac{\partial L(\theta)}{\partial z^{l+1}} \\ (W_{\cdot,2}^{l+1})^\mathsf{T} \frac{\partial L(\theta)}{\partial z^{l+1}} \\ \dots \\ (W_{\cdot,n}^{l+1})^\mathsf{T} \frac{\partial L(\theta)}{\partial z^{l+1}} \end{pmatrix}$$

$$\frac{\partial L(\theta)}{\partial a_j^l} = (W^{l+1})^\mathsf{T} \frac{\partial L(\theta)}{\partial z^{l+1}}$$
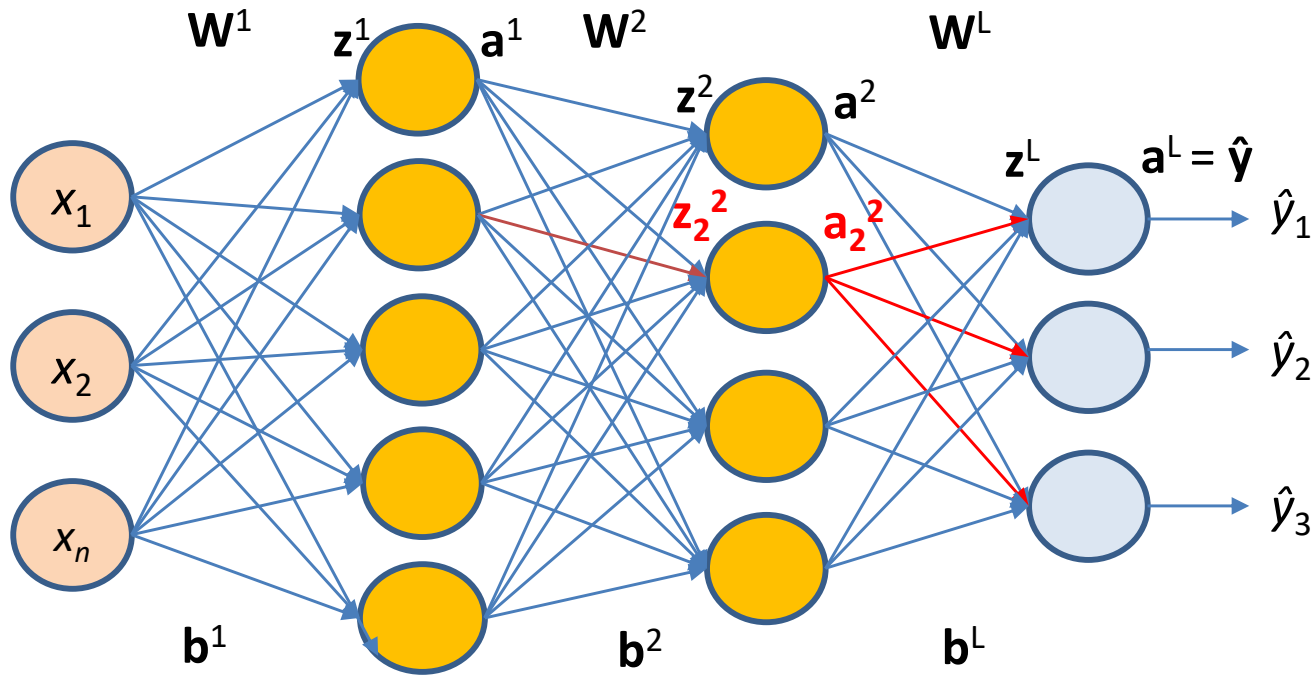
How to compute $\frac{\partial L(\theta)}{\partial z^{l+1}}$ for $l \leq L\text{-}1$

# Gradients w.r.t. hidden units

$$\frac{\partial L(\theta)}{\partial z^l} = \begin{bmatrix} \dfrac{\partial L(\theta)}{\partial z_1^l} \\ \dfrac{\partial L(\theta)}{\partial z_2^l} \\ \dots \\ \dfrac{\partial L(\theta)}{\partial z_n^l} \end{bmatrix}$$

$$\frac{\partial L(\theta)}{\partial z_j^l} = \frac{\partial L(\theta)}{\partial a_j^l} * \frac{\partial a_j^l}{\partial z_j^l} = \frac{\partial L(\theta)}{\partial a_j^l} * \frac{\partial g(z_j^l)}{\partial z_j^l}$$

$$\frac{\partial L(\theta)}{\partial z^l} = \begin{bmatrix} \dfrac{\partial L(\theta)}{\partial a_1^l} * \dfrac{\partial g(z_1^l)}{\partial z_1^l} \\ \dfrac{\partial L(\theta)}{\partial a_2^l} * \dfrac{\partial g(z_2^l)}{\partial z_2^l} \\ \dots \\ \dfrac{\partial L(\theta)}{\partial a_n^l} * \dfrac{\partial g(z_n^l)}{\partial z_n^l} \end{bmatrix} = \frac{\partial L(\theta)}{\partial a_j^l} \odot \left[ \dots, \dots, \frac{\partial g(z_j^l)}{\partial z_j^l}, \dots \right]$$

# Gradients w.r.t. Parameters



$$\frac{\partial L(\theta)}{\partial w^l} = \frac{\partial L(\theta)}{\partial z^l} * \frac{\partial z^l}{\partial w^l}$$

$$z^l = W^l \, a^{l-1} + b^l$$
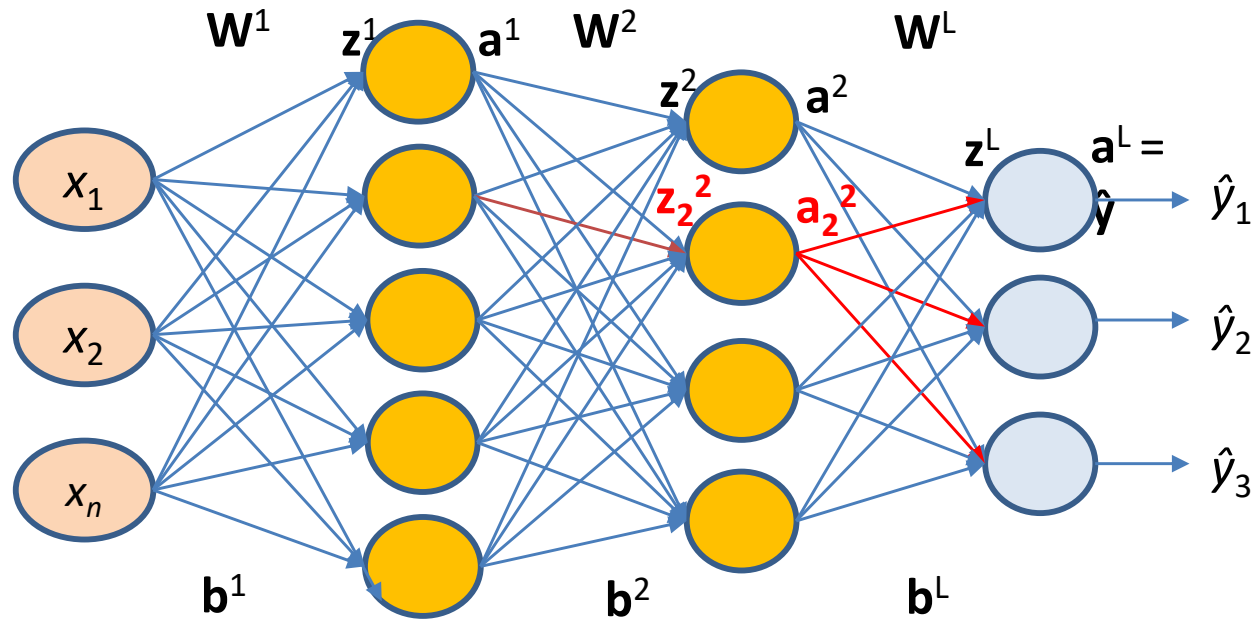
# Gradients w.r.t. Parameters



$$\frac{\partial L(\theta)}{\partial w^l} = \frac{\partial L(\theta)}{\partial z^l} * \frac{\partial z^l}{\partial w^l}$$

$$z^l = W^l \, a^{l-1} + b^l$$

$$\begin{bmatrix} z_1^l \\ z_2^l \\ .... \\ z_n^l \end{bmatrix} = \begin{bmatrix} w_{11}^l \; w_{12}^l \; ... \; w_{1n}^l \\ w_{21}^l \; w_{22}^l \; ... \; w_{2n}^l \\ ... \\ w_{n1}^l \; w_{n2}^l \; ... \; w_{nn}^l \end{bmatrix} \begin{bmatrix} a_1^{l-1} \\ a_2^{l-1} \\ \\ a_n^{l-1} \end{bmatrix} + \begin{bmatrix} b_1^{l-1} \\ b_2^{l-1} \\ \\ b_n^{l-1} \end{bmatrix}$$

# Gradients w.r.t. Parameters

$$\frac{\partial L(\theta)}{\partial w^l} = \frac{\partial L(\theta)}{\partial z^l} * \frac{\partial z^l}{\partial w^l}$$

$$\begin{bmatrix} z_1^l \\ z_2^l \\ .... \\ z_n^l \end{bmatrix} = \begin{bmatrix} w_{11}^l\ w_{12}^l\ ...\ w_{1n}^l \\ w_{21}^l\ w_{22}^l\ ...\ w_{2n}^l \\ ... \\ w_{n1}^l\ w_{n2}^l\ ...\ w_{nn}^l \end{bmatrix} \begin{bmatrix} a_1^{l-1} \\ a_2^{l-1} \\ \\ a_n^{l-1} \end{bmatrix} + \begin{bmatrix} b_1^{l-1} \\ b_2^{l-1} \\ \\ b_n^{l-1} \end{bmatrix}$$

$$\frac{\partial z^l}{\partial w_{12}^l} = a_2^{l-1} \qquad\qquad \frac{\partial z^l}{\partial w_{pq}^l} = a_q^{l-1}$$

$$\frac{\partial L(\theta)}{\partial w_{pq}^l} = \frac{\partial L(\theta)}{\partial z^l} * a_q^{l-1}$$

# Gradients w.r.t. Parameters

$$\frac{\partial L(\theta)}{\partial w_{pq}^l} = \frac{\partial L(\theta)}{\partial z^l} * a_q^{l-1}$$

$$\frac{\partial L(\theta)}{\partial w^l} = \begin{pmatrix} \dfrac{\partial L(\theta)}{\partial w_{00}^l} & \dfrac{\partial L(\theta)}{\partial w_{01}^l} & \cdots & \dfrac{\partial L(\theta)}{\partial w_{0n-1}^l} \\ \cdots & \cdots & \cdots & .. \\ \cdots & \cdots & & \dfrac{\partial L(\theta)}{\partial w_{n-1\,n-1}^l} \end{pmatrix}$$

# Backpropagation Equations

$$\frac{\partial L}{\partial z^l} = (w^{l+1})^T \cdot \frac{\partial L^l}{\partial z^{l+1}} * \sigma'(z^l)$$

$$\frac{\partial L}{\partial w^l} = \frac{\partial L^l}{\partial z^l} \cdot (a^{l-1})^T$$

$$\frac{\partial L}{\partial b^l} = \frac{\partial L^l}{\partial z^l}$$

\* - element-wise multiplication

. - dot product

*l* – layer number

(assuming sigmoid activations)

# Pseudo code: Gradient Descent

$t$:=0;

$max\_iterations$:=1000;

Initialize $\theta_0 := [W^1_0, \dots W^L_0, b^1_0 \dots b^L_0]$;

**while** $t{+}{+} < max\_iterations$ **do**

$\quad a^1, a^2 \dots a^{L-1}, z^1, z^2 \dots z^L, \hat{y} = forward\_propagation(\theta_t)$;

$\quad \nabla\theta_t = back\_propagation(a^1, a^2 \dots a^{L-1}, z^1, z^2 \dots z^L, y, \hat{y})$;

$\quad \theta_{t+1} := \theta_t - \eta\nabla\theta_t$;

**end**

# Pseudo code: Forward Propagation

**for** $v = 0$ to $L$-1 **do**

$\quad\quad z^v = b^v + W^v a^{v-1};$

$\quad\quad a^v = g(z^v);$

**end**;

$z^L = b^L + W^L a^{L-1};$

$\hat{y} = O(z^L);$

Do a forward propagation and compute all $a^i$'s, $z^i$'s, and $\hat{y}$

# Pseudo code: Back Propagation

//Compute output gradient:

$\nabla_z^L L(\theta) = -(e(y) - \hat{y})$;

**for** *v = L* to 1 **do**

   //*Compute gradients w.r.t. parameters*

   $\nabla_W^v L(\theta) = \nabla_z^v L(\theta) \, a^{v-1}$;

   $\nabla_b^v L(\theta) = \nabla_z^v L(\theta)$ ;

   //*Compute gradients w.r.t. layer below*

   $\nabla_a^{v-1} L(\theta) = W^v \nabla_z^v L(\theta)$;

   //*Compute gradients w.r.t. layer below (pre-activation)*

   $\nabla_z^{v-1} L(\theta) = \nabla_a^{v-1} L(\theta) \odot [\ldots\ldots g'(z^{v-1,j})\ldots\ldots]$;

**end**