# Lecture 14: Backpropagation for training Neural Nets

- A neural network is a parametric model

- We will tune its parameters using optimization techniques such as SGD, ADAM, etc.

- To find suitable values for the parameters $\underline{\Theta}$,

$$\underline{\Theta} = \begin{bmatrix} \text{vec}(\underline{\underline{W}}^{(1)}) \\ \underline{b}^{(1)} \\ \vdots \\ \text{vec}(\underline{\underline{W}}^{(L)}) \\ \underline{b}^{(L)} \end{bmatrix}$$

we will solve the optimization problem:

$$\hat{\underline{\Theta}} = \underset{\underline{\Theta}}{\arg\min} \; \underline{J(\underline{\Theta})}$$

$$\text{cost function}$$

where $J(\underline{\Theta}) = \frac{1}{N} \sum_{i=1}^{N} \underline{L(y_i, f(\underline{x}_i; \underline{\Theta}))}$

loss function
for data point $(\underline{x}_i, y_i)$

$$J(\underline{\Theta}) = \frac{1}{N} \sum_{i=1}^{N} \underbrace{L(y_i, f(\underline{x}_i; \underline{\Theta}))}_{}$$

loss function
for data point $(\underline{x}_i, y_i)$

– The functional form of the loss function depends on the problem

- Regression problems typically use squared error loss

$$L(y, f(\underline{x}; \underline{\Theta})) = \left(y - \underline{f(\underline{x}; \underline{\Theta})}\right)^2$$

output of a neural net

- Multi-class classification problems typically use cross-entropy loss

(with M classes)

softmax

$$L(y, f(\underline{x}; \underline{\Theta})) = -\ln g_m\left(\underline{f}(\underline{x}; \underline{\Theta})\right)$$

$\underline{f}$ is a vector of $M \times 1$

$$= -\ln g_y\left(\underline{f}(\underline{x}; \underline{\Theta})\right)$$

we use training data label $y$
as an index variable to select the correct logit

$$\hat{\underline{\theta}} = \arg\min_{\underline{\theta}} J(\underline{\theta}) \quad \text{where} \quad J(\underline{\theta}) = \frac{1}{N} \sum_{i=1}^{N} L\left(y_i, f(\underline{x}_i; \underline{\theta})\right)$$

- These optimization problems cannot be solved in closed form

  $\rightarrow$ Numerical optimization algorithms have to be used

- Numerical optimization update parameters in an iterative manner

  In deep learning, one typically uses gradient-descent algorithms

Step 1 : Pick an initial guess $\underline{\theta}^{(0)}$

Step 2 : Calculate gradient of the cost function w.r.t $\underline{\theta}^{(t)}$, $t = 0, 1, 2, ...$

Step 3 : Update the parameters as $\underline{\theta}^{(t+1)} \leftarrow \underline{\theta}^{(t)} - \gamma \nabla_{\underline{\theta}} J(\underline{\theta}^{(t)})$

Step 4 : Terminate when some criterion is fulfilled, and take the last $\underline{\theta}^{(t)}$ as $\hat{\underline{\theta}}$

# Computational Challenges

1> **Large datasets (N very large)**

    — The number of datapoints N is large in deep learning applications

    — Makes computation of the cost function gradient very costly, due to the sum

    — We resort to using a <u>random subset of data</u> to update parameters

                                ↳ minibatch gradient descent

2> **Large number of parameters $\underline{\theta}$**

    — The dimension of the parameter vector $\underline{\theta}$ is very large in deep learning

    — To efficiently calculate the gradient $\nabla_{\underline{\theta}} J(\underline{\theta}^{(t)})$,

      we need to apply chain rule of calculus

this will be done using the Backpropagation algorithm

# Univariate chain rule

$$z = wx + b \quad \text{← parameters}$$

$$\hat{y} = \sigma(z) \quad \text{← non-linear activation function}$$

— Let's compute the derivative of loss function w.r.t parameters $w$ and $b$

$$L = (y - \hat{y})^2 \quad \text{(e.g. sigmoid)}$$

loss function

$$L = (y - \sigma(wx+b))^2$$

$$\frac{\partial L}{\partial w} = \frac{\partial}{\partial w}(y - \sigma(wx+b))^2$$

$$= (y - \sigma(wx+b)) \frac{\partial}{\partial w}(y - \sigma(wx+b))$$

$$= -(y - \sigma(wx+b)) \, \sigma'(wx+b) \frac{\partial}{\partial w}(wx+b)$$

$$= -(y - \sigma(wx+b)) \, \sigma'(wx+b) \, x$$

$$\frac{\partial L}{\partial b} = \frac{\partial}{\partial b}(y - \sigma(wx+b))^2$$

$$= (y - \sigma(wx+b)) \frac{\partial}{\partial b}(y - \sigma(wx+b))$$

$$= -(y - \sigma(wx+b)) \, \sigma'(wx+b)$$

## Disadvantages of this approach

- Calculations are very cumbersome

  A lot of terms have been copied from one line to the next

- Final expression has repeated terms

# Univariate chain rule

A more structured approach of chain rule would be:

1) Compute the loss

$$z = wx + b$$

$$\hat{y} = \sigma(z)$$

$$L = (y - \hat{y})^2$$

2) Compute the derivatives

$$\frac{\partial L}{\partial \hat{y}} = -2(y - \hat{y})$$

computed from previous step

$$\frac{\partial L}{\partial z} = \frac{\partial L}{\partial \hat{y}} \frac{d\hat{y}}{dz} = \frac{\partial L}{\partial \hat{y}} \sigma'(z)$$

evaluated at current step

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial w} = \frac{\partial L}{\partial z} x$$

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

This form of computation is clean and has no repeated expressions!

# Computational graph

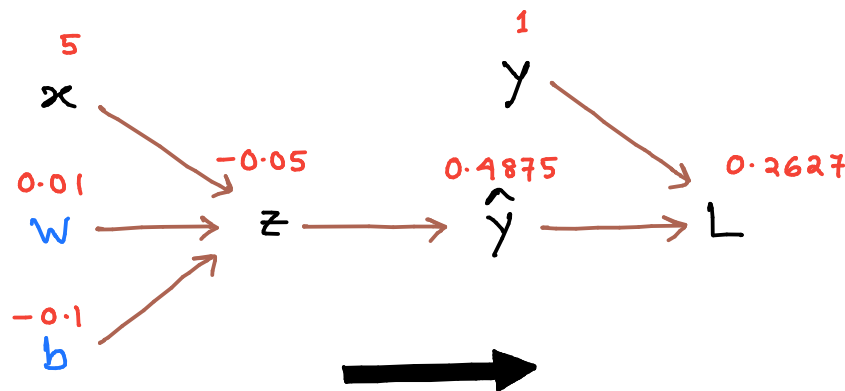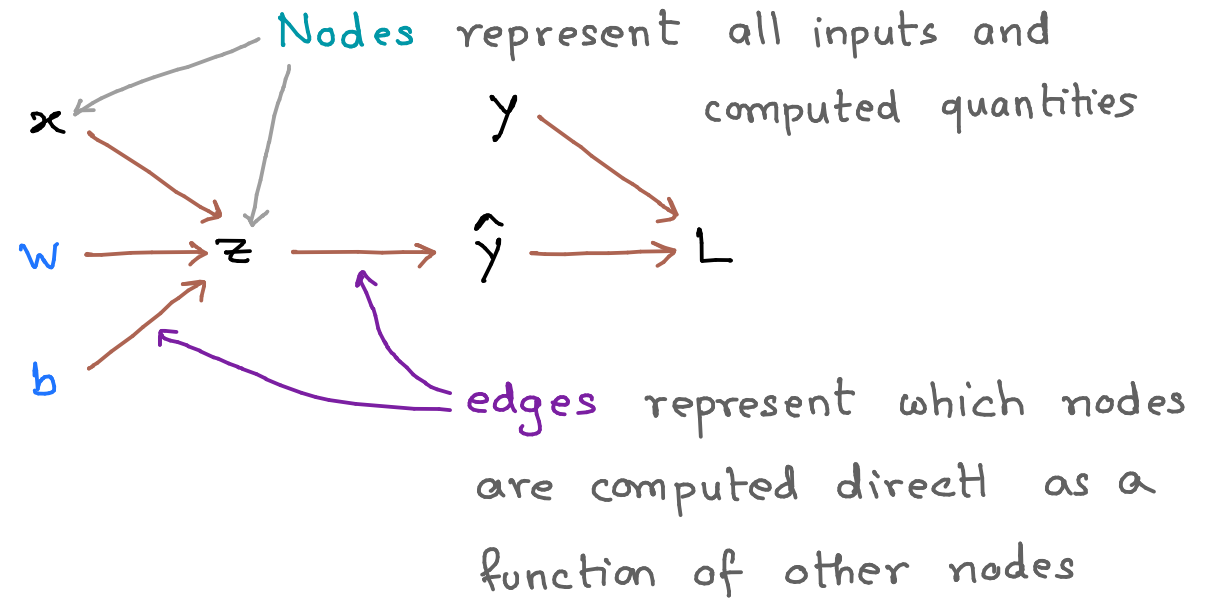- The computations can be plotted using a *computational graph*

**1)** Compute the loss

$$z = wx + b$$

$$\hat{y} = \sigma(z)$$
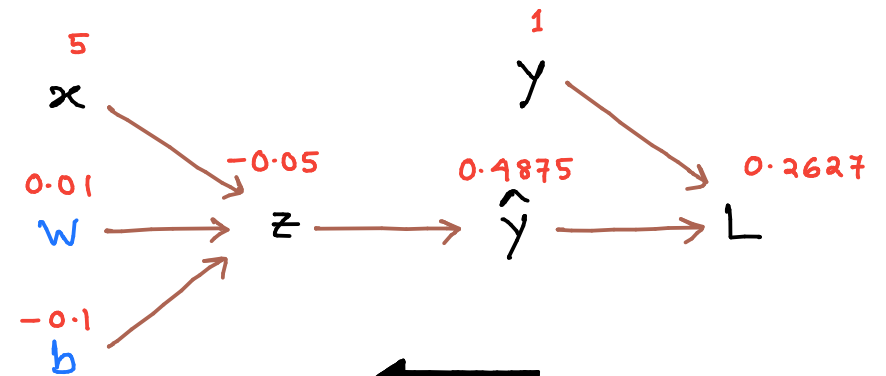
$$L = (y - \hat{y})^2$$

sigmoid activation
$$\frac{e^z}{1 + e^z}$$

Nodes represent all inputs and computed quantities

x → z ← w, b → z → ŷ → L ← y

edges represent which nodes are computed directl as a function of other nodes

5
x
0.01
w
−0.1
b
−0.05
z
0.4875
ŷ
1
y
0.2627
L

Forward pass for loss:   Parents ⟶ Children

For backprop, we will use notation

$$\bar{v} = \frac{\partial L}{\partial v} \quad , \quad \bar{x} = \frac{\partial L}{\partial x} \quad , \text{ etc.}$$

$x$ — 5

$y$ — 1

$w$ — 0.01

$b$ — $-0.1$

$z$ — $-0.05$

$\hat{y}$ — 0.4875

$L$ — 0.2627

Backward pass for gradients

Parents ⟵ Children

1) Compute the loss

$$z = wx + b$$
$$\hat{y} = \sigma(z)$$
$$L = (y - \hat{y})^2$$

2) Compute the derivatives

$$\frac{\partial L}{\partial \hat{y}} = -2(y - \hat{y})$$

$$\frac{\partial L}{\partial z} = \frac{\partial L}{\partial \hat{y}} \frac{d\hat{y}}{dz} = \frac{\partial L}{\partial \hat{y}} \sigma'(z)$$

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial w} = \frac{\partial L}{\partial z} x$$

$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial z} \frac{\partial z}{\partial b} = \frac{\partial L}{\partial z}$$
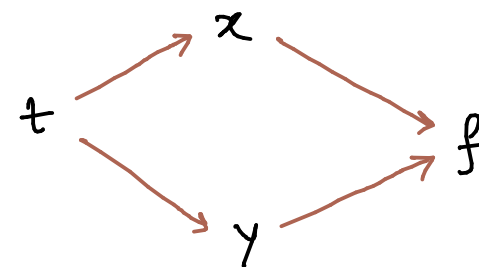
$$\bar{\hat{y}} = -2(y - \hat{y})$$

$$\bar{z} = \bar{\hat{y}} \, \sigma'(z)$$

$$\bar{w} = \bar{z} \, x$$

$$\bar{b} = \bar{z}$$

# Multivariate Chain Rule
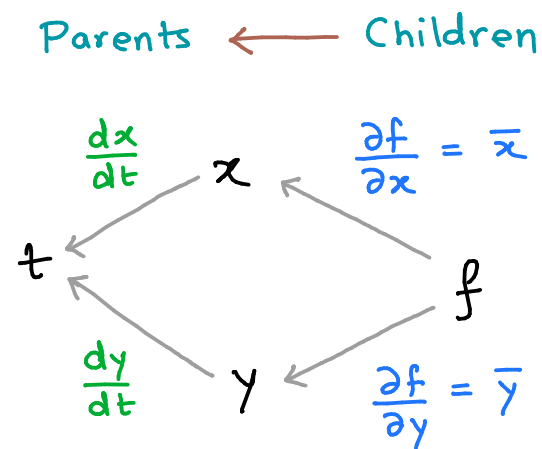
- Suppose we have a function $f(x(t), y(t))$

$$\frac{\partial f}{\partial t} = \frac{\partial f}{\partial x} \frac{dx}{dt} + \frac{\partial f}{\partial y} \frac{dy}{dt}$$

- In the context of gradient computation (backward pass)

these values will be computed first

$$\frac{\partial f}{\partial t} = \frac{\partial f}{\partial x} \frac{dx}{dt} + \frac{\partial f}{\partial y} \frac{dy}{dt}$$

Parents $\longleftarrow$ Children

$$\frac{\partial f}{\partial t} = \frac{\partial f}{\partial x} \frac{dx}{dt} + \frac{\partial f}{\partial y} \frac{dy}{dt}$$

These will be evaluated next
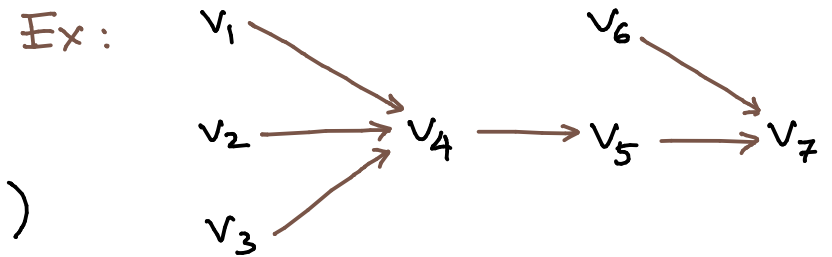
In our notation

$$\bar{t} = \bar{x} \frac{dx}{dt} + \bar{y} \frac{dy}{dt}$$

$\frac{dx}{dt}$  $x$  $\frac{\partial f}{\partial x} = \bar{x}$

$t$  $f$

$\frac{dy}{dt}$  $y$  $\frac{\partial f}{\partial y} = \bar{y}$

# Backpropagation Algorithm

- Let $v_1, v_2, \ldots, v_P$ be a topological ordering of the computational graph (i.e. where parents come before children)

Ex:



- $v_P$ denotes the variable we are trying to compute the derivatives of

In our case $v_P \equiv L$ (loss function)

**Forward pass**

Compute values

For $i = 1, \ldots, P$

     Compute $v_i$ as a function of Parents $(v_i)$
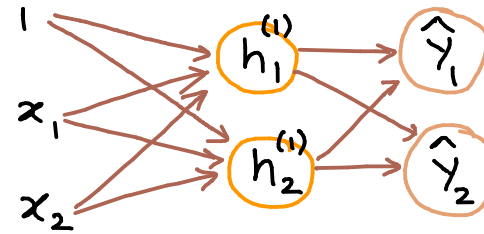
**Backward pass**

Compute derivatives

Treat $\overline{v_P} = \dfrac{\partial v_P}{\partial v_P} = 1$

For $i = P-1, \ldots, 1$

$$\overline{v_i} = \sum_{j \in \text{Children}(v_j)} \overline{v_j} \, \frac{\partial v_j}{\partial v_i}$$

# Backpropagation for Neural Nets

Neural net with 1-hidden layer
(with multiple outputs)



## Forward pass
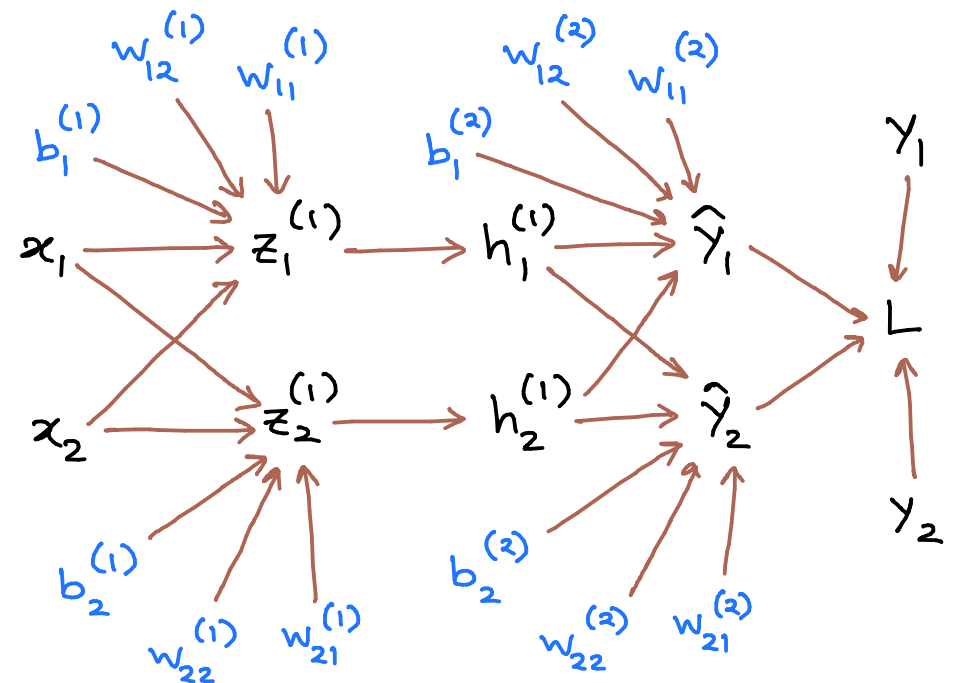(to compute loss)

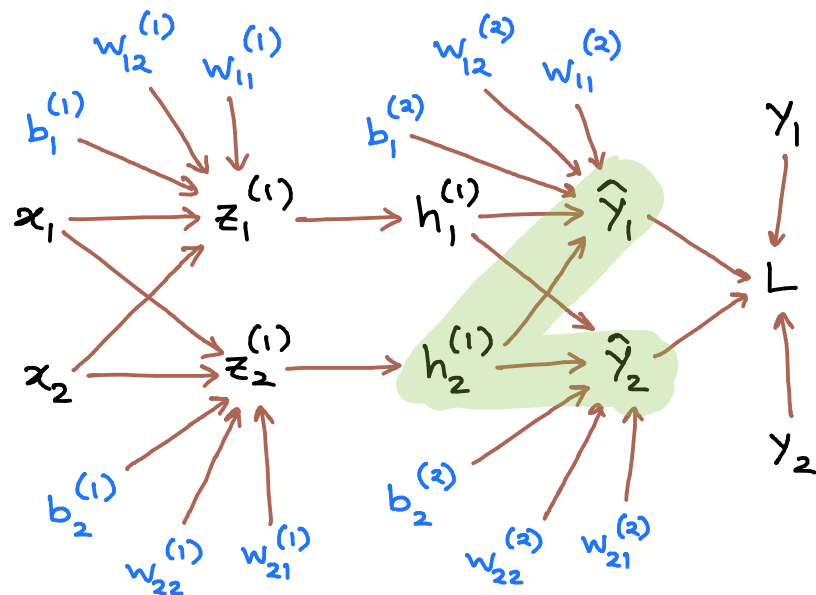$$z_i^{(1)} = \sum_j w_{ij}^{(1)} x_j + b_i^{(1)}$$

$$h_i^{(1)} = \sigma\left(z_i^{(1)}\right)$$

$$\hat{y}_k = \sum_i w_{ki}^{(2)} h_i^{(1)} + b_k^{(2)}$$

$$L = \sum_k \left(y_k - \hat{y}_k\right)^2$$

## Computational graph

**Forward pass** (to compute loss)

$$z_i^{(1)} = \sum_j w_{ij}^{(1)} x_j + b_i^{(1)}$$

$$h_i^{(1)} = \sigma\left(z_i^{(1)}\right)$$

$$\hat{y}_k = \sum_i w_{ki}^{(2)} h_i^{(1)} + b_k^{(2)}$$

$$L = \sum_k \left(y_k - \hat{y}_k\right)^2$$

$$\frac{\partial L}{\partial h_2^{(1)}} = \frac{\partial L}{\partial \hat{y}_1}\frac{\partial \hat{y}_1}{\partial h_2^{(1)}} + \frac{\partial L}{\partial \hat{y}_2}\frac{\partial \hat{y}_2}{\partial h_2^{(1)}}$$

**Backward pass** (to compute gradients)

$$\bar{L} = 1$$

$$\bar{\hat{y}}_k = \bar{L}\,\frac{\partial L}{\partial \hat{y}_k} = -2\bar{L}\left(y_k - \hat{y}_k\right)$$

$$\bar{w}_{ki}^{(2)} = \bar{\hat{y}}_k\,\frac{\partial \hat{y}_k}{\partial w_{ki}^{(2)}} = \bar{\hat{y}}_k\, h_i^{(1)}$$

$$\bar{b}_k^{(2)} = \bar{\hat{y}}_k\,\frac{\partial \hat{y}_k}{\partial b_k^{(2)}} = \bar{\hat{y}}_k$$

$$\bar{h}_i^{(1)} = \sum_k \bar{\hat{y}}_k\,\frac{\partial \hat{y}_k}{\partial h_i^{(1)}} = \sum_k \bar{\hat{y}}_k\, w_{ki}^{(2)}$$

$$\bar{z}_i^{(1)} = \bar{h}_i^{(1)}\,\frac{\partial h_i^{(1)}}{\partial z_i^{(1)}} = \bar{h}_i^{(1)}\,\sigma'\left(z_i^{(1)}\right)$$

$$\bar{w}_{ij}^{(1)} = \bar{z}_i^{(1)}\,\frac{\partial z_i^{(1)}}{\partial w_{ij}^{(1)}} = \bar{z}_i^{(1)}\, x_j$$

$$\bar{b}_i^{(1)} = \bar{z}_i^{(1)}\,\frac{\partial z_i^{(1)}}{\partial b_i^{(1)}} = \bar{z}_i^{(1)}$$

# Vectorized form of BackProp

- Computational graphs showing individual units are cumbersome

- Instead draw graphs over the vectorized variables

$$\underline{\underline{W}}^{(1)} \qquad \underline{\underline{W}}^{(2)} \qquad y$$

$$\underline{x} \longrightarrow \underline{z}^{(1)} \longrightarrow \underline{h}^{(1)} \longrightarrow \hat{\underline{y}} \longrightarrow L$$

$$\underline{b}^{(1)} \qquad \underline{b}^{(2)}$$

## Backprop rules

### Forward Pass

Compute values

For $i = 1, \dots, P$

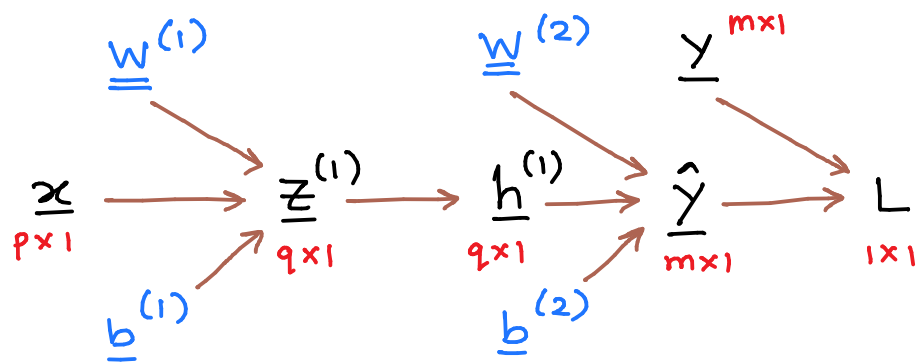    Compute $\underline{v}_i$ as a function of Parents $(\underline{v}_i)$

### Backward Pass: Compute derivatives

Treat $\overline{\underline{v}}_P = 1$

For $i = P-1, \dots, 1$

$$\overline{\underline{v}}_i = \sum_{j \in \text{Children}(\underline{v}_j)} \overline{\underline{v}}_j \frac{\partial \underline{v}_j}{\partial \underline{v}_i}$$

## Forward Pass

Network diagram (top left):

$\underline{\underline{W}}^{(1)}$, $\underline{x}$ ($p \times 1$), $\underline{b}^{(1)}$ $\to$ $\underline{z}^{(1)}$ ($q \times 1$) $\to$ $\underline{h}^{(1)}$ ($q \times 1$), $\underline{\underline{W}}^{(2)}$, $\underline{b}^{(2)}$ $\to$ $\hat{\underline{y}}$ ($m \times 1$), $\underline{y}^{m \times 1}$ $\to$ $L$ ($1 \times 1$)

$$\underset{q \times 1}{\underline{z}^{(1)}} = \underset{q \times p}{\underline{\underline{W}}^{(1)}} \ \underset{p \times 1}{\underline{x}} + \underset{q \times 1}{\underline{b}^{(1)}}$$

$$\underset{q \times 1}{\underline{h}^{(1)}} = \sigma(\underset{q \times 1}{\underline{z}^{(1)}})$$

$$\underset{m \times 1}{\hat{\underline{y}}} = \underset{m \times q}{\underline{\underline{W}}^{(2)}} \ \underset{q \times 1}{\underline{h}^{(1)}} + \underset{m \times 1}{\underline{b}^{(2)}}$$

$$\underset{1 \times 1}{L} = \underset{1 \times m}{\left(\underline{y} - \hat{\underline{y}}\right)^T} \ \underset{m \times 1}{\left(\underline{y} - \hat{\underline{y}}\right)}$$

## Backward Pass

$$\underset{1 \times 1}{\bar{L}} = 1$$

$$\frac{\partial L}{\partial \hat{y}} \rightarrow \underset{m \times 1}{\bar{\hat{\underline{y}}}} = -2\bar{L}\underset{m \times 1}{\left(\underline{y} - \hat{\underline{y}}\right)}$$

$$\frac{\partial L}{\partial \underline{\underline{W}}^{(2)}} \rightarrow \underset{m \times q}{\bar{\underline{\underline{W}}}^{(2)}} = \underset{m \times 1}{\bar{\hat{\underline{y}}}} \ \underset{1 \times q}{\underline{h}^{(1)^T}}$$

$$\underset{m \times 1}{\bar{\underline{b}}^{(2)}} = \underset{m \times 1}{\bar{\hat{\underline{y}}}}$$

$$\frac{\partial L}{\partial \underline{h}^{(1)}} \rightarrow \underset{q \times 1}{\bar{\underline{h}}^{(1)}} = \underset{q \times m}{\bar{\underline{\underline{W}}}^{(2)^T}} \ \underset{m \times 1}{\bar{\hat{\underline{y}}}}$$

$$\underset{q \times 1}{\bar{\underline{z}}^{(1)}} = \underset{q \times 1}{\bar{\underline{h}}^{(1)}} \circ \underset{q \times 1}{\sigma'(\underline{z}^{(1)})}$$

elementwise product

$$\frac{\partial L}{\partial \underline{\underline{W}}^{(1)}} \rightarrow \underset{q \times p}{\bar{\underline{\underline{W}}}^{(1)}} = \underset{q \times 1}{\bar{\underline{z}}^{(1)}} \ \underset{1 \times p}{\underline{x}^T}$$

$$\underset{q \times 1}{\bar{\underline{b}}^{(1)}} = \underset{q \times 1}{\bar{\underline{z}}^{(1)}}$$

# Back Prop example in vectorized form



$$w^{(1)} \quad w^{(2)} \quad y$$
$$x \longrightarrow z^{(1)} \longrightarrow h^{(1)} \longrightarrow \hat{y} \longrightarrow L$$
$$b^{(1)} \quad b^{(2)}$$

## Forward Pass

$$z^{(1)} = w^{(1)} x + b^{(1)}$$

$$h^{(1)} = \sigma(z^{(1)})$$

$$\hat{y} = w^{(2)} h^{(1)} + b^{(2)}$$

$$L = (y - \hat{y})^T (y - \hat{y})$$

## Backward Pass

$$\overline{L} = 1$$

$$\overline{\hat{y}} = -2\overline{L}(y - \hat{y})$$

$$\overline{W}^{(2)} = \overline{\hat{y}} \, h^{(1)^T}$$

$$\overline{b}^{(2)} = \overline{\hat{y}}$$

$$\overline{h}^{(1)} = \overline{W}^{(2)^T} \hat{y}$$

$$\overline{z}^{(1)} = \overline{h}^{(1)} \circ \sigma'(z^{(1)})$$

$$\overline{W}^{(1)} = \overline{z}^{(1)} x^T$$

$$\overline{b}^{(1)} = \overline{z}^{(1)}$$

- Backprop in neural networks are commonly implemented as matrix-vector multiplications

- These matrix-vector multiplications are called vector Jacobian products (VJPs)

# Closing Remarks

- Backprop is based on the computational graph, and it basically works backwards through the graph, applying the chain rule at each node

- Backprop is used to train most neural nets you will find these days

- Even optimization algorithms much fancier than gradient descent (such as second-order methods) use backprop to compute gradients

- Once the derivatives w.r.t. the weights and biases are computed using backprop, the updates are applied to the weights and biases using some optimization scheme

$$\underline{\underline{W}}^{(t+1)} \leftarrow \underline{\underline{W}}^{(t)} - \eta \left.\frac{\partial J}{\partial \underline{\underline{W}}}\right|_{\underline{\underline{W}}^{(t)}} \qquad \underline{b}^{(t+1)} \leftarrow \underline{b}^{(t)} - \eta \left.\frac{\partial J}{\partial \underline{b}}\right|_{\underline{b}^{(t)}}$$

- Hand-calculation of derivatives are replaced with automatic differentiation