

APL 745

Backpropagation

Dr. Rajdip Nayek

Block V, Room 418D

Department of Applied Mechanics

Indian Institute of Technology Delhi

E-mail: rajdipn@am.iitd.ac.in

Overview

- We've seen that multilayer neural networks are powerful. But **how to actually learn (or optimize) the weights and biases?**

Update
formula

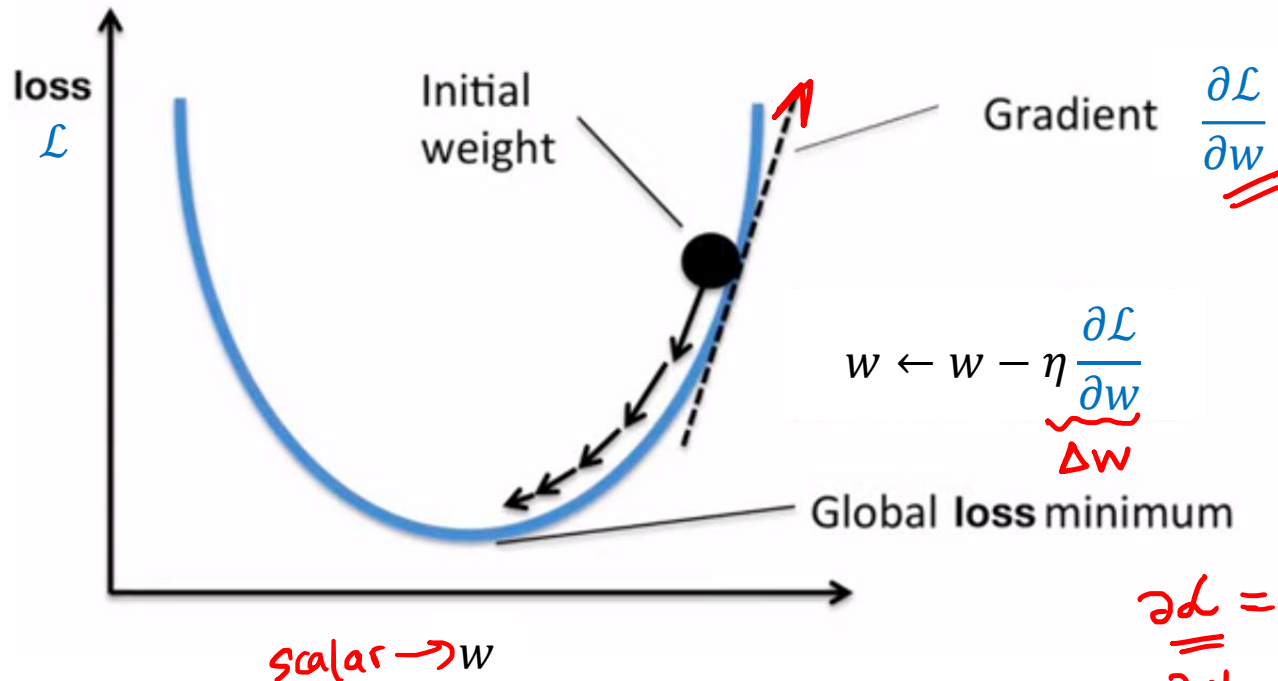
$$\mathbf{w} \leftarrow \mathbf{w} - \eta \frac{\partial \mathcal{L}}{\partial \mathbf{w}}$$

- **Backpropagation** is the central algorithm that is **used for computing gradients**
 - Backpropagation is an efficient use of the **Chain Rule for derivatives**
 - It's an instance of **reverse mode automatic differentiation**

N.B. Lecture slides mostly follow the material of CSC421 by Roger Grosse and Jimmy Ba

Recap: Gradient descent

- Recall that gradient descent moves opposite to the direction of gradient



$$\frac{\partial \mathcal{L}}{\partial w} = \frac{1}{N_{\text{train}}} \sum_{i=1}^{N_{\text{train}}} (t^{(i)} - y^{(i)})^2$$

- We want to compute the gradient of the loss function $\partial \mathcal{L} / \partial w$, which is the vector of partial derivatives and is calculated **by averaging over all training examples**
- In this lecture, we will focus on computing the gradients of loss function for a single training example

Univariate Chain Rule

- We already know of the univariate chain rule
- Recall if $f(x)$ and $x(t)$ are univariate functions, then

$$\frac{d}{dt}f(x(t)) = \frac{df}{dx} \frac{dx}{dt} \quad \leftarrow \text{chain rule}$$

- Consider a univariate logistic least-squares model

Supervised

Preactivation $z = wx + b$ \leftarrow input

Model pred $y = \sigma(z)$ \leftarrow output

loss $\ell = \frac{1}{2}(t - y)^2$ \leftarrow single training example

(x, t)

- Let's compute the loss derivatives

Univariate Chain Rule

- Let's compute the loss derivatives w.r.t. w and b using calculus

$$\ell = \frac{1}{2} (t - \sigma(wx + b))^2$$

$$\frac{\partial \ell}{\partial w} = \frac{\partial}{\partial w} \left[\frac{1}{2} (t - \sigma(wx + b))^2 \right]$$

$$= \frac{1}{2} \frac{\partial}{\partial w} (t - \sigma(wx + b))^2$$

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

$$= -(t - \sigma(wx + b)) \sigma'(wx + b) \frac{\partial}{\partial w} (wx + b)$$

$$= -(t - \sigma(\underline{wx + b})) \sigma'(\underline{wx + b}) x$$

$$\frac{\partial \ell}{\partial b} = \frac{\partial}{\partial b} \left[\frac{1}{2} (t - \sigma(wx + b))^2 \right]$$

$$= \frac{1}{2} \frac{\partial}{\partial b} (t - \sigma(wx + b))^2$$

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

$$= -(t - \sigma(wx + b)) \sigma'(wx + b) \frac{\partial}{\partial b} (wx + b)$$

$$= -(t - \sigma(wx + b)) \sigma'(wx + b)$$

What are the disadvantages of this approach?

- The calculations are very cumbersome. In this derivation, we had to copy lots of terms from one line to the next
- The final expressions have lots of repeated terms

Univariate Chain Rule

A more structured way of doing it

1) Compute the loss:

$$z = wx + b$$

$$y = \sigma(z)$$

$$\rightarrow \ell = \frac{1}{2}(t - \underline{y})^2$$

This form of computation is clean!

- No repeated expressions

2) Compute the derivatives:

$$\begin{aligned} \frac{\partial \ell}{\partial y} &= -(t - y) \\ \frac{\partial \ell}{\partial z} &= \frac{\partial \ell}{\partial y} \frac{dy}{dz} = \frac{\partial \ell}{\partial y} \sigma'(z) \\ \frac{\partial \ell}{\partial w} &= \frac{\partial \ell}{\partial z} \frac{\partial z}{\partial w} = \frac{\partial \ell}{\partial z} x \\ \frac{\partial \ell}{\partial b} &= \frac{\partial \ell}{\partial z} \frac{\partial z}{\partial b} = \frac{\partial \ell}{\partial z} \end{aligned}$$

obtained from previous step

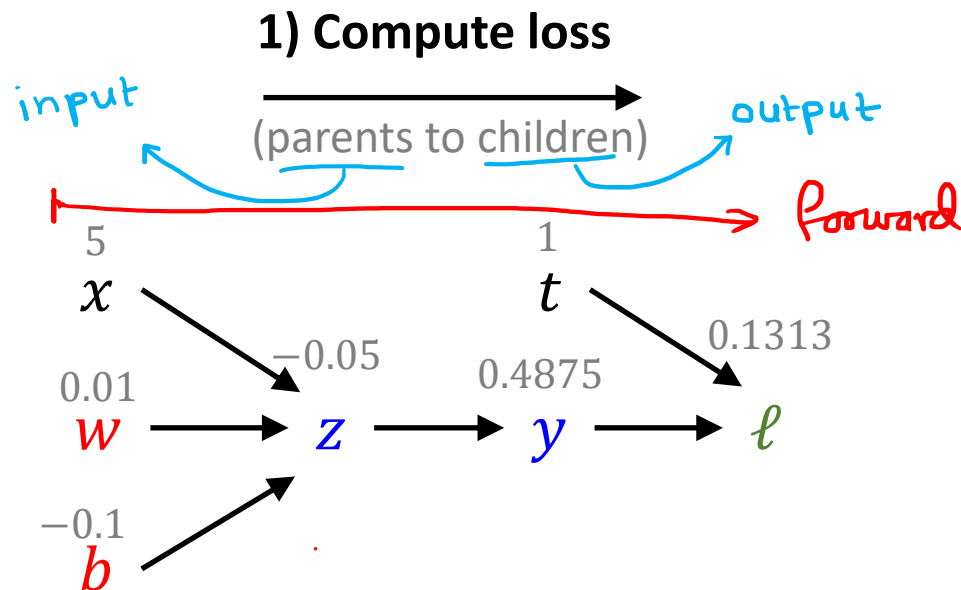
evaluated at current step

Univariate Chain Rule

- We can plot these computations using a **computation graph**
- The **nodes** represent all the inputs and computed quantities
- The **edges** represent which nodes are computed directly as a function of other nodes

1) Compute loss

$$\begin{aligned} \rightarrow z &= wx + b \\ \rightarrow y &= \sigma(z) \\ \rightarrow \ell &= \frac{1}{2}(t - y)^2 \end{aligned}$$



Univariate Chain Rule

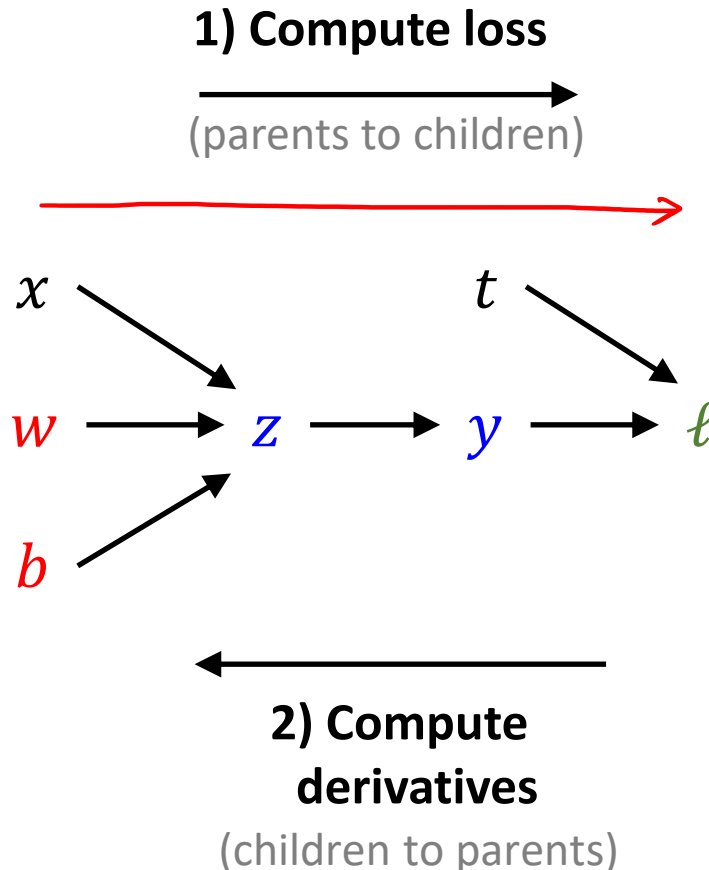
- We can plot these computations using a **computation graph**
- The **nodes** represent all the inputs and computed quantities
- The **edges** represent which nodes are computed directly as a function of other nodes

1) Compute loss

$$z = wx + b$$

$$y = \sigma(z)$$

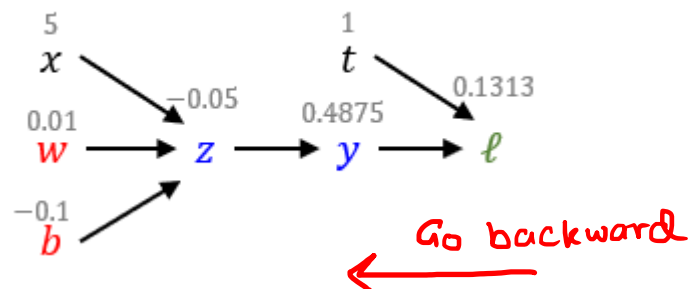
$$\ell = \frac{1}{2}(t - y)^2$$



Univariate Chain Rule

A weird but slightly convenient notation

- Usual notation: $\nabla_v \ell = \frac{\partial \ell}{\partial v}$
- Instead, use $\bar{v} = \nabla_v \ell = \frac{\partial \ell}{\partial v}$



2) Compute the derivatives:

1) Compute the loss:

$$z = wx + b$$

$$y = \sigma(z)$$

$$\ell = \frac{1}{2}(t - y)^2$$

$$\frac{\partial \ell}{\partial y} = -(t - y)$$

$$\frac{\partial \ell}{\partial z} = \frac{\partial \ell}{\partial y} \frac{dy}{dz} = \frac{\partial \ell}{\partial y} \sigma'(z)$$

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

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

$$\bar{y} = -(t - y)$$

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

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

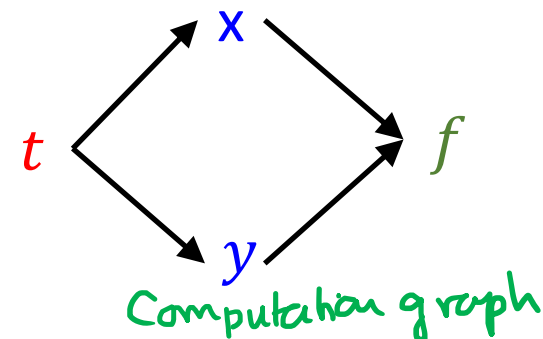
$$\bar{b} = \bar{z}$$

Note: \bar{w} used here should not be confused with augmented \bar{w} used in linear regression/classification

Multivariate Chain Rule

- Suppose we have a function $f(x(t), y(t))$. Then

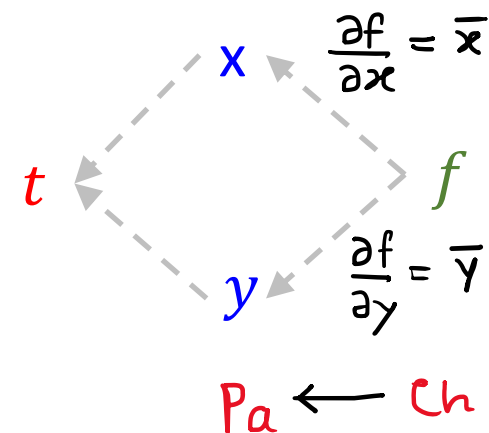
$$\frac{d}{dt} f(x(t), y(t)) = \frac{\partial f}{\partial x} \frac{dx}{dt} + \frac{\partial f}{\partial y} \frac{dy}{dt} \quad \checkmark$$



- In the context of backpropagation (backward pass)

$$\frac{d}{dt} f(x(t), y(t)) = \frac{\partial f}{\partial x} \frac{dx}{dt} + \frac{\partial f}{\partial y} \frac{dy}{dt}$$

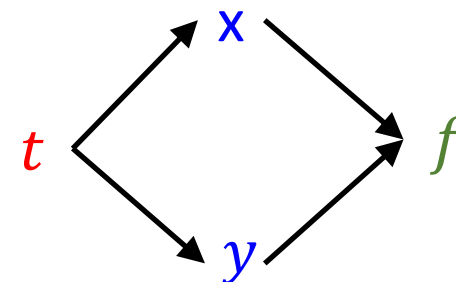
Values computed first



Multivariate Chain Rule

- Suppose we have a function $f(x(t), y(t))$. Then

$$\frac{d}{dt} f(x(t), y(t)) = \frac{\partial f}{\partial x} \frac{dx}{dt} + \frac{\partial f}{\partial y} \frac{dy}{dt}$$



- In the context of backpropagation (backward pass)

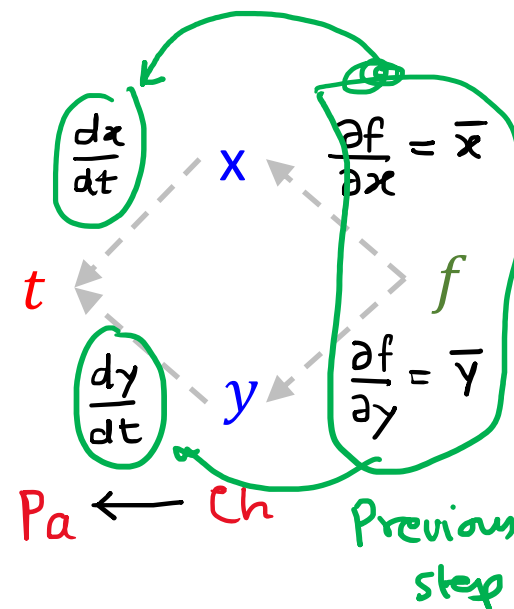
Mathematical expressions
to be evaluated

$$\frac{d}{dt} f(x(t), y(t)) = \frac{\partial f}{\partial x} \frac{dx}{dt} + \frac{\partial f}{\partial y} \frac{dy}{dt}$$

Values already computed

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

- In our notation:

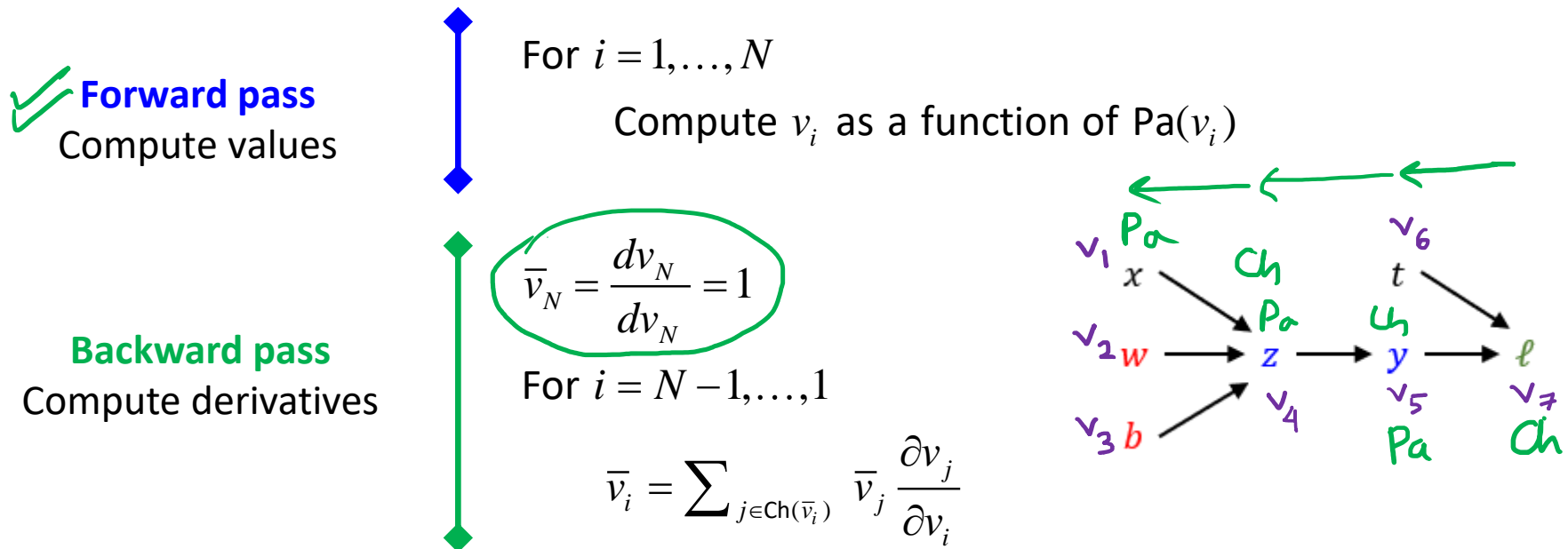


Backpropagation

Full backpropagation algorithm:

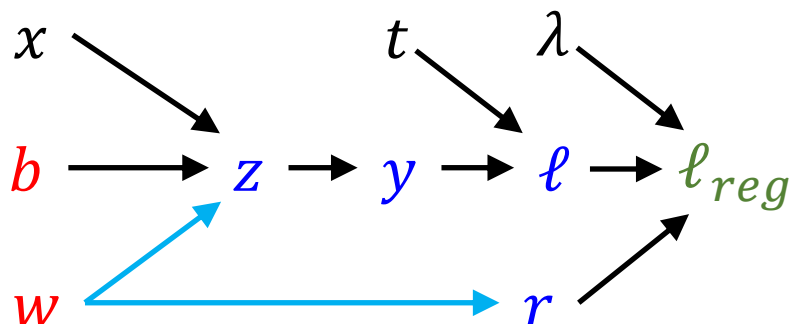
Let v_1, \dots, v_N be a topological ordering of the computation graph
(i.e. where parents come before children)

v_N denotes the variable we are trying to compute derivatives of (e.g. loss function)



Backpropagation

Example: Logistic least-squares regression



Forward pass

Compute values

$$z = wx + b$$

$$y = \sigma(z)$$

$$\ell = \frac{1}{2}(t - y)^2$$

$$r = \frac{1}{2}w^2$$

$$\ell_{reg} = \ell + \lambda r$$

Backward pass

Compute derivatives

$$\bar{v}_i = \frac{\partial \ell_{reg}}{\partial v_i}$$

$$\bar{\ell}_{reg} = 1$$

$$\bar{r} = \bar{\ell}_{reg} \frac{\partial \ell_{reg}}{\partial r}$$

$$= \bar{\ell}_{reg} \lambda$$

$$\frac{\partial \ell_{reg}}{\partial \ell} = \bar{\ell} = \bar{\ell}_{reg} \frac{\partial \ell_{reg}}{\partial \ell}$$

$$= \bar{\ell}_{reg}$$

$$\frac{\partial \ell_{reg}}{\partial y} = \bar{\ell} \frac{\partial \ell}{\partial y}$$

$$= -\bar{\ell} (t - y)$$

$$\bar{z} = \bar{y} \frac{dy}{dz}$$

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

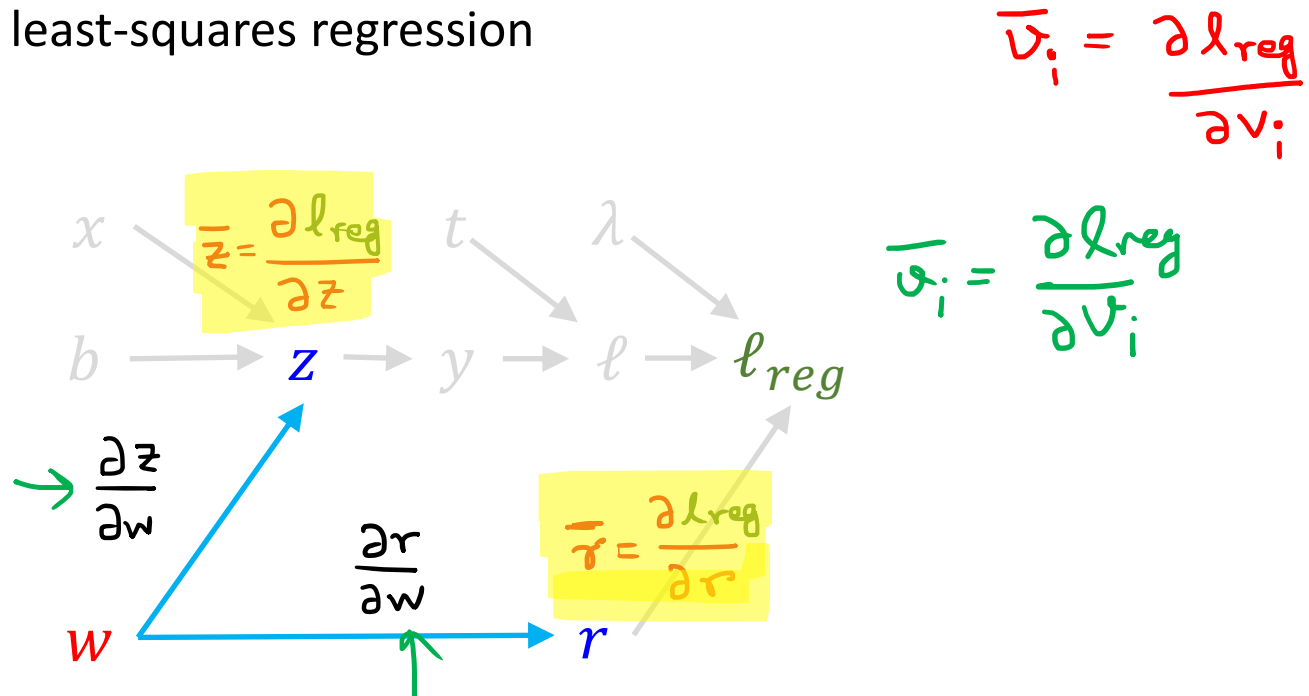
$$\bar{b} = \bar{z} \frac{\partial z}{\partial b} = \bar{z}$$

$$\bar{w} = \bar{z} \frac{\partial z}{\partial w} + \bar{r} \frac{\partial r}{\partial w}$$

$$= \bar{z}x + \bar{r}w$$

Backpropagation

Example: Logistic least-squares regression



$$z = wx + b$$

$$y = \sigma(z)$$

$$\ell = \frac{1}{2}(t - y)^2$$

$$r = \frac{1}{2}w^2$$

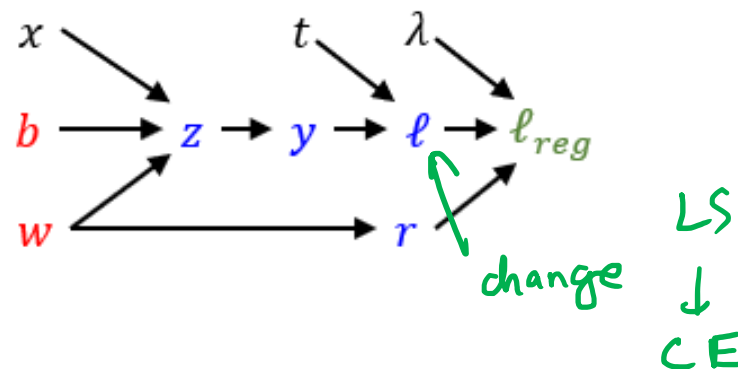
$$\ell_{reg} = \ell + \lambda r$$

$$\begin{aligned} \frac{\partial \ell_{reg}}{\partial w} &= \bar{w} = \bar{z} \frac{\partial z}{\partial w} + \bar{r} \frac{\partial r}{\partial w} \\ &= \bar{z}x + \bar{r}w \end{aligned}$$

Backpropagation

Example: Logistic least-squares regression

Derivatives via backprop



- The derivation, and the final result, are **much cleaner** and **efficient**. There are no redundant computations here.
- The procedure is **modular**: it is broken down into small chunks that can be reused for other computations. For instance, if we want to change the loss function, we'd only have to modify the formula for \bar{y}

$$\bar{l}_{reg} = 1$$

$$\begin{aligned}\bar{r} &= \bar{l}_{reg} \frac{\partial l_{reg}}{\partial r} \\ &= \bar{l}_{reg} \lambda\end{aligned}$$

$$\begin{aligned}\bar{l} &= \bar{l}_{reg} \frac{\partial l_{reg}}{\partial l} \\ &= \bar{l}_{reg}\end{aligned}$$

$$\begin{aligned}\bar{y} &= \bar{l} \frac{\partial l}{\partial y} \\ &= -\bar{l} (t - y)\end{aligned}$$

$\frac{\partial l_{reg}}{\partial t}$

$$\begin{aligned}\bar{z} &= \bar{y} \frac{dy}{dz} \\ &= \bar{y} \sigma'(z)\end{aligned}$$

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

$$\begin{aligned}\bar{w} &= \bar{z} \frac{\partial z}{\partial w} + \bar{r} \frac{\partial r}{\partial w} \\ &= \bar{z}x + \bar{r}w\end{aligned}$$

$$\bar{\lambda} = \frac{\partial l_{reg}}{\partial \lambda}$$

Backpropagation for MLP

Multilayer Perceptron (multiple outputs)

One hidden layer

Want to compute gradients of loss function wrt all weights and biases

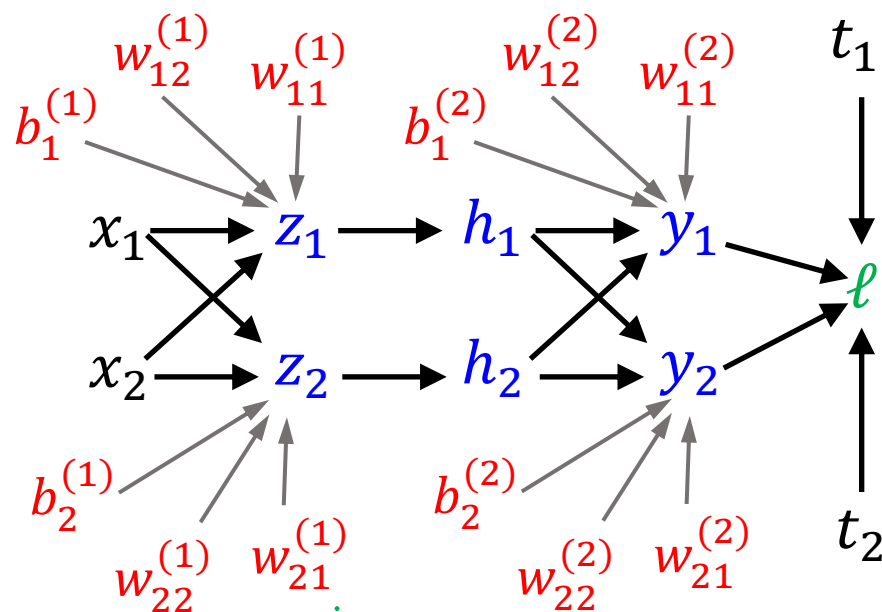
Forward pass

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

$$h_i = \sigma(z_i)$$

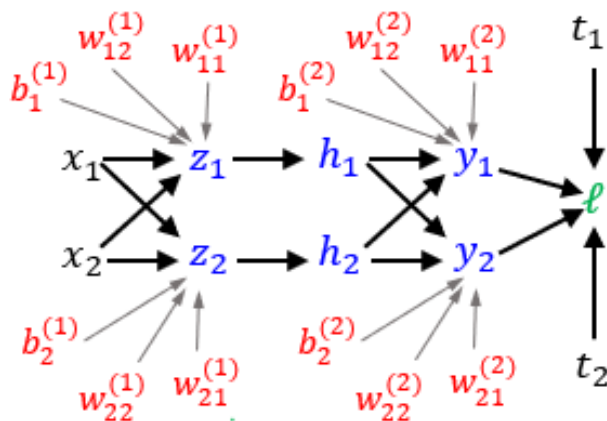
$$y_k = \sum_i w_{ki}^{(2)} h_i + b_k^{(2)}$$

$$\ell = \frac{1}{2} \sum_k (t_k - y_k)^2$$



Backpropagation for MLP

Multilayer Perceptron (multiple outputs):



Forward pass

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

$$h_i = \sigma(z_i)$$

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

$$\ell = \frac{1}{2} \sum_k (t_k - y_k)^2 \quad \checkmark$$

Backward pass

$$\bar{\ell} = 1$$

$$\bar{y}_k = \bar{\ell} \frac{\partial \ell}{\partial y_k} = -\bar{\ell} (t_k - y_k) \quad \checkmark$$

$$\bar{w}_{ki}^{(2)} = \bar{y}_k \frac{\partial y_k}{\partial w_{ki}^{(2)}} = \bar{y}_k h_i \quad \checkmark$$

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

$$\rightarrow \bar{h}_i = \sum_k \bar{y}_k \frac{\partial y_k}{\partial h_i} = \sum_k \bar{y}_k w_{ki}^{(2)} \quad \checkmark$$

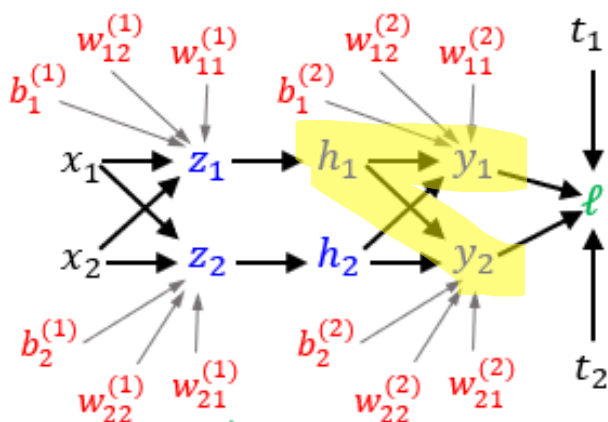
$$\bar{z}_i = \bar{h}_i \frac{\partial h_i}{\partial z_i} = \bar{h}_i \sigma'(z_i) \quad \checkmark$$

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

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

Backpropagation for MLP

Multilayer Perceptron (multiple outputs):



Forward pass

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

$$h_i = \sigma(z_i)$$

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

$$\ell = \frac{1}{2} \sum_k (t_k - y_k)^2$$

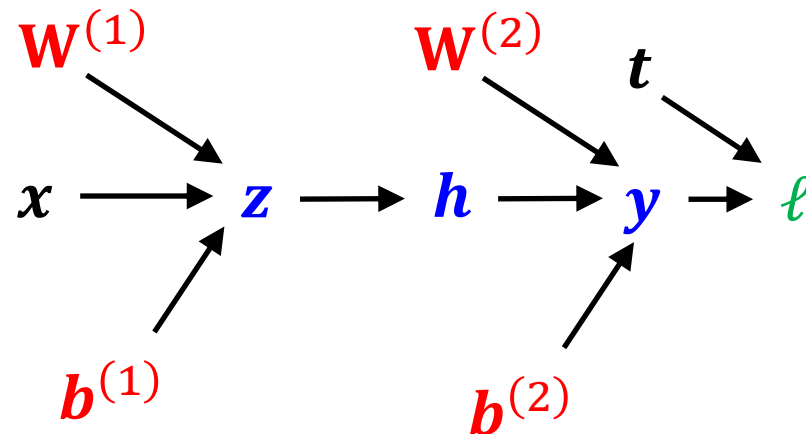
Full break

$$\bar{h}_1 = \bar{y}_1 \frac{\partial y_1}{\partial h_1} + \bar{y}_2 \frac{\partial y_2}{\partial h_1} = \bar{y}_1 w_{11}^{(2)} + \bar{y}_2 w_{21}^{(2)}$$

Red arrows indicate the flow of gradients: from $\frac{\partial \ell}{\partial h_1}$ to \bar{h}_1 , from $\frac{\partial \ell}{\partial y_1}$ to \bar{y}_1 , and from $\frac{\partial \ell}{\partial y_2}$ to \bar{y}_2 . Red downward arrows point from the text "Full break" to the terms \bar{y}_1 and \bar{y}_2 in the equation.

Vector form

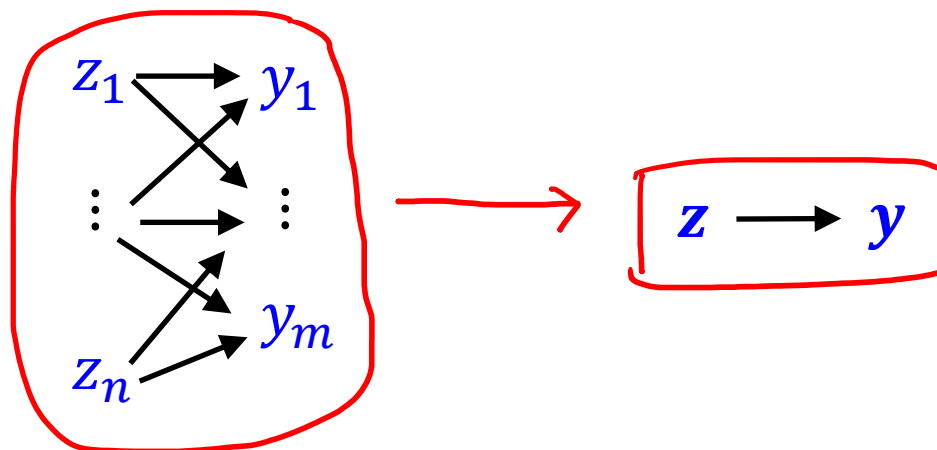
- Computation graphs showing individual units are cumbersome
- We can draw graphs over the **vectorized** variables



- We pass the gradients back in the same way as for the scalar-valued nodes

Vector form

- Consider this partial computation graph:



- Backprop rules: $\underline{\bar{z}}_j = \sum_{k=1}^m \bar{y}_k \frac{\partial y_k}{\partial z_j}$

where $\frac{\partial y}{\partial z}$ is the Jacobian matrix

$$\frac{\partial y}{\partial z} = \begin{pmatrix} \frac{\partial y_1}{\partial z_1} & \dots & \frac{\partial y_1}{\partial z_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial y_m}{\partial z_1} & \dots & \frac{\partial y_m}{\partial z_n} \end{pmatrix}$$

$$\bar{z} = \left(\frac{\partial y}{\partial z} \right)^T \bar{y}$$

$\bar{y} = \frac{\partial L}{\partial y}$ $m \times 1$
 $n \times 1$ $n \times m$ $m \times 1$
matrix vector
 $m \times n$ VJP

Vector form

Full backpropagation algorithm (vector form):

Let $\mathbf{v}_1, \dots, \mathbf{v}_N$ be a topological ordering of the computation graph
(i.e. parents come before children)

\mathbf{v}_N denotes the variable we are trying to compute derivatives of (e.g. loss function)

Forward pass

Compute values

For $i = 1, \dots, N$

Compute \mathbf{v}_i as a function of $\text{Pa}(\mathbf{v}_i)$

Backward pass

Compute derivatives

$\bar{\mathbf{v}}_N = 1$

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

$$\bar{\mathbf{v}}_i = \sum_{j \in \text{Ch}(\bar{\mathbf{v}}_i)} \left(\frac{\partial \mathbf{v}_j}{\partial \mathbf{v}_i} \right)^T \bar{\mathbf{v}}_j$$

Vector form

MLP example in vectorized form:

Backward
pass

$$\bar{\ell} = 1$$

$$1 \times 1$$

$$\bar{y} = -\bar{\ell} (t - y)$$

$$m \times 1 \quad m \times 1$$

$$\bar{W}^{(2)} = \bar{y} h^T$$

$$m \times n \quad 1 \times n$$

$$\bar{b}^{(2)} = \bar{y}$$

$$m \times 1 \quad m \times 1$$

$$\bar{h} = W^{(2)T} \bar{y}$$

$$n \times 1 \quad n \times m \quad m \times 1$$

$$\bar{z} = \bar{h} \circ \sigma'(z)$$

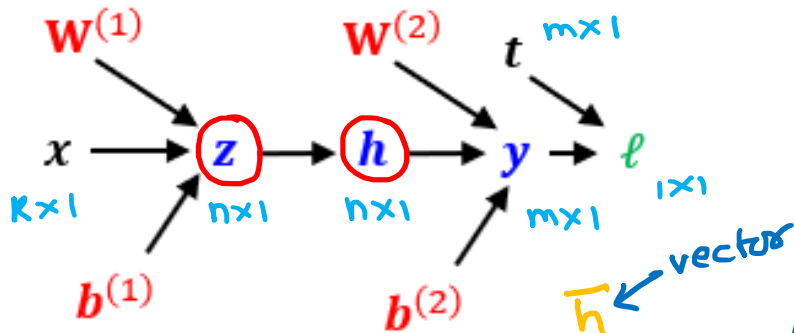
$$n \times 1 \quad n \times 1 \quad n \times 1$$

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

$$n \times k \quad 1 \times k$$

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

$$n \times 1 \quad n \times 1$$



$$\bar{z} = \frac{\partial \ell}{\partial z} = \left[\frac{\partial \ell}{\partial h} \right] \left[\frac{\partial h}{\partial z} \right]$$

vector

vector

Forward pass

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

$$n \times 1 \quad n \times k \quad k \times 1 \quad n \times 1$$

$$h = \sigma(z)$$

$$n \times 1 \quad n \times 1$$

$$y = W^{(2)} h + b^{(2)}$$

$$m \times 1 \quad m \times n \quad n \times 1 \quad m \times 1$$

$$\ell = \frac{1}{2} (t - y)^T (t - y)$$

$$1 \times m \quad m \times 1 \quad m \times 1$$

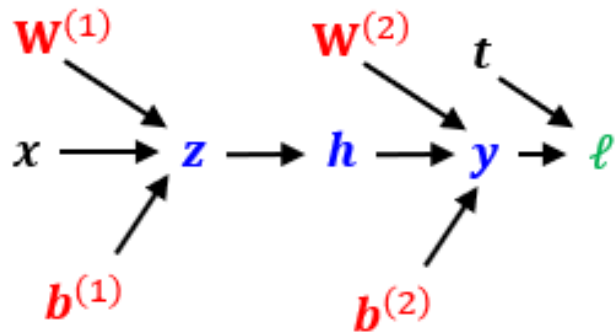
Elementwise
multiplication

$$\frac{\partial \ell}{\partial W^{(1)}} = \frac{\partial \ell}{\partial y} \times \left[\frac{\partial y}{\partial W^{(1)}} \right]$$

0

Vector form

MLP example in vectorized form:



Forward pass

$$z = \mathbf{W}^{(1)} x + \mathbf{b}^{(1)}$$

$$h = \sigma(z)$$

$$y = \mathbf{W}^{(2)} h + \mathbf{b}^{(2)}$$

$$\ell = \frac{1}{2} (t - y)^T (t - y)$$

Backward
pass

$$\bar{\ell} = 1$$

$$\frac{\partial \ell}{\partial y} \rightarrow \bar{y} = -\bar{\ell} (t - y) \quad \frac{\partial \ell}{\partial y}$$

$$\bar{\mathbf{W}}^{(2)} = \bar{y} h^T$$

$$\bar{\mathbf{b}}^{(2)} = \bar{y}$$

$$\bar{h} = \mathbf{W}^{(2)T} \bar{y}$$

Elementwise multiplication

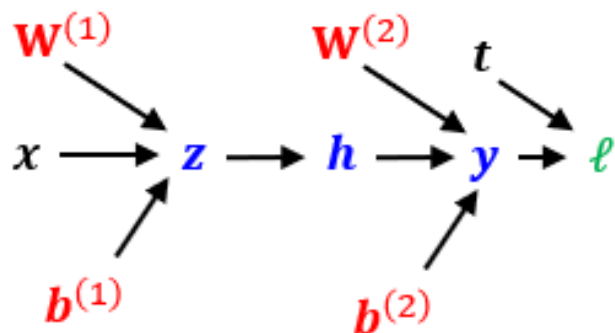
$$\bar{z} = \bar{h} \circ \sigma'(z)$$

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

$$\bar{\mathbf{b}}^{(1)} = \bar{z}$$

Backpropagation in MLPs

MLP example in vectorized form:



Forward pass

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

$$h = \sigma(z)$$

$$y = W^{(2)}h + b^{(2)}$$

$$\ell = \frac{1}{2}(t - y)^T(t - y)$$

Backward pass

$$\bar{\ell} = 1$$

$$\bar{y} = -\bar{\ell}(t - y)$$

$$\bar{W}^{(2)} = \bar{y}h^T$$

$$\bar{b}^{(2)} = \bar{y}$$

$$\bar{h} = W^{(2)T}\bar{y}$$

$$\bar{z} = \bar{h} \circ \sigma'(z)$$

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

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

- Backpropagation in MLPs are commonly implemented as matrix-vector multiplications
- Here, these matrix-vector multiplications can be called vector Jacobian products (VJPs)

Closing remarks

- Backprop is used to train most of the neural nets you will find today
- Even optimization algorithms much fancier than gradient descent (e.g. second-order methods) use backprop to compute the gradients
- Backprop is based on the computation graph, and it basically works backwards through the graph, applying the chain rule at each node
- 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

$$\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial \mathcal{L}}{\partial \mathbf{W}}$$
$$\mathbf{b} \leftarrow \mathbf{b} - \eta \frac{\partial \mathcal{L}}{\partial \mathbf{b}}$$

- However, here we wrote out the computation graph and calculated the derivatives by hand ourselves; this hand calculation is avoided by automatic differentiation