

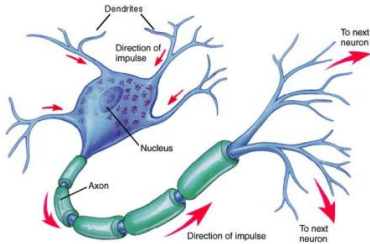
# Topic 10: NEURAL NETWORKS

---

STAT 37710/CAAM 37710/CMSC 35400 Machine Learning  
Risi Kondor, The University of Chicago

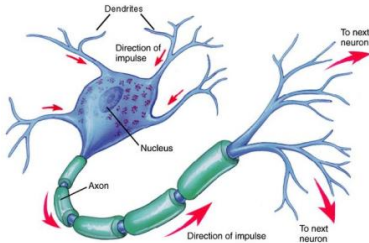


# Intro to neural networks



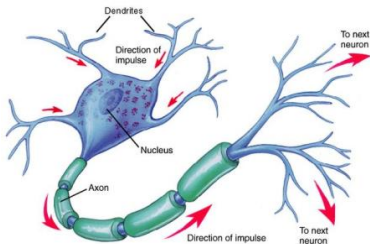
- The human brain has  $\sim 10^{11}$  neurons, each connected to  $\sim 10^4$  others.

# Intro to neural networks



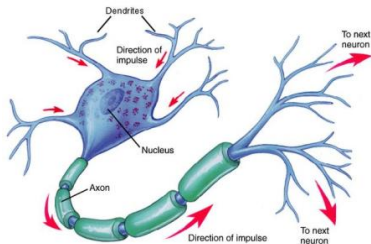
- The human brain has  $\sim 10^{11}$  neurons, each connected to  $\sim 10^4$  others.
- Inputs come from the dendrites, are aggregated in the soma, if the neuron starts firing impulses propagated to other neurons via axons.

# Intro to neural networks



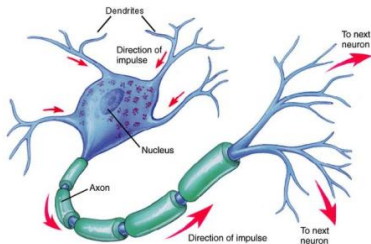
- The human brain has  $\sim 10^{11}$  neurons, each connected to  $\sim 10^4$  others.
  - Inputs come from the dendrites, are aggregated in the soma, if the neuron starts firing impulses propagated to other neurons via axons.
- 
- Neurons learn by changing the connection strengths of their synapses.

# Intro to neural networks



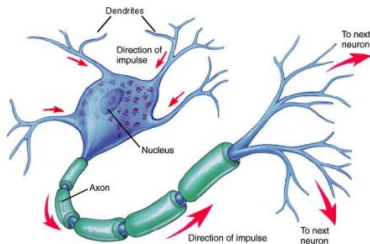
- The human brain has  $\sim 10^{11}$  neurons, each connected to  $\sim 10^4$  others.
- Inputs come from the dendrites, are aggregated in the soma, if the neuron starts firing impulses propagated to other neurons via axons.
- Neurons learn by changing the connection strengths of their synapses.
- Information storage in the nervous system is “distributed”.

# Intro to neural networks



- The human brain has  $\sim 10^{11}$  neurons, each connected to  $\sim 10^4$  others.
  - Inputs come from the dendrites, are aggregated in the soma, if the neuron starts firing impulses propagated to other neurons via axons.
- 
- Neurons learn by changing the connection strengths of their synapses.
  - Information storage in the nervous system is “distributed”.
  - The response time of the brain is quite fast, so the “depth” of the network can’t be very great. (clear layer by layer organization in the visual system).

# Intro to neural networks



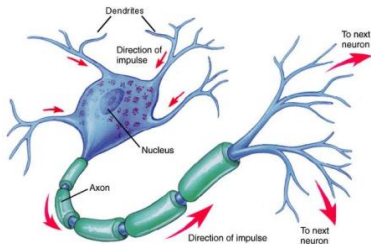
- The human brain has  $\sim 10^{11}$  neurons, each connected to  $\sim 10^4$  others.
- Inputs come from the dendrites, are aggregated in the soma, if the neuron starts firing impulses propagated to other neurons via axons.

- Neurons learn by changing the connection strengths of their synapses.
- Information storage in the nervous system is “distributed”.
- The response time of the brain is quite fast, so the “depth” of the network can’t be very great. (clear layer by layer organization in the visual system).

**IDEA:** Humans seem to be okay at learning, so why not try to replicate this in a computer?



# Intro to neural networks

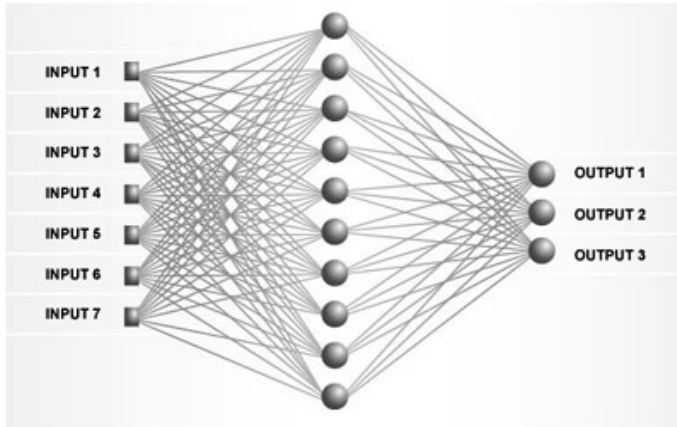


- The human brain has  $\sim 10^{11}$  neurons, each connected to  $\sim 10^4$  others.
- Inputs come from the dendrites, are aggregated in the soma, if the neuron starts firing impulses propagated to other neurons via axons.

- Neurons learn by changing the connection strengths of their synapses.
- Information storage in the nervous system is “distributed”.
- The response time of the brain is quite fast, so the “depth” of the network can’t be very great. (clear layer by layer organization in the visual system).

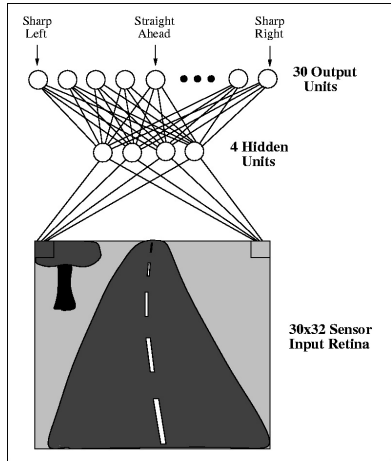
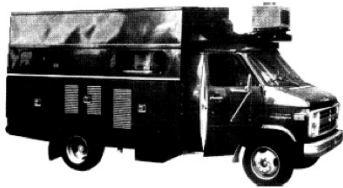
**IDEA:** Humans seem to be okay at learning, so why not try to replicate this in a computer? Goes back to the early days of AI, many successes and failures.

# Multilayer artificial neural net



Question: But what should the individual neurons do and how should they learn?

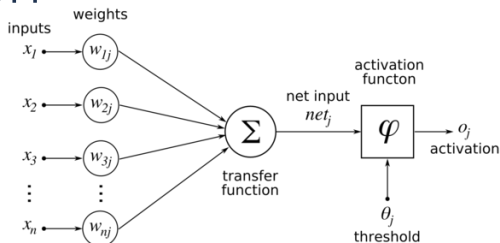
# A success: ALVINN [Pomerleau '95]



Drove unassisted from Pittsburgh to NYC on the highway.

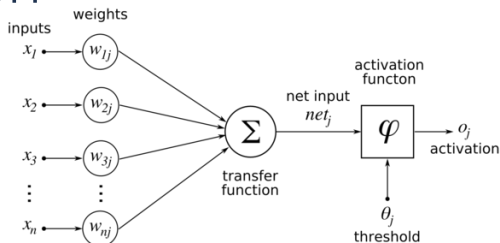
# A model neuron

$$o = \sigma\left(\theta + \sum_p w_p x_p\right),$$



# A model neuron

$$o = \sigma\left(\theta + \sum_p w_p x_p\right),$$

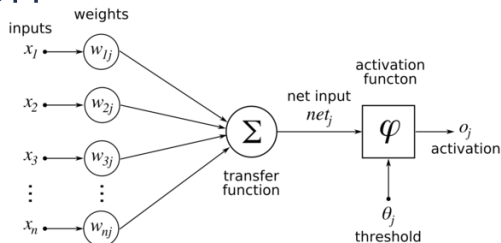


## Notation:

- $x_i$ : the  $i$ 'th input

# A model neuron

$$o = \sigma\left(\theta + \sum_p w_p x_p\right),$$

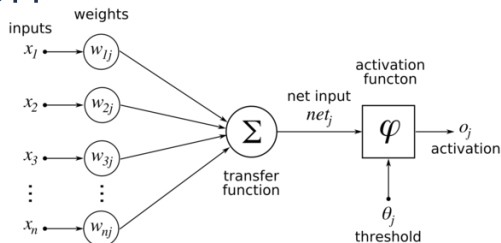


## Notation:

- $x_i$ : the  $i$ 'th input
- $w_i$ : the corresponding synaptic weight

# A model neuron

$$o = \sigma\left(\theta + \sum_p w_p x_p\right),$$

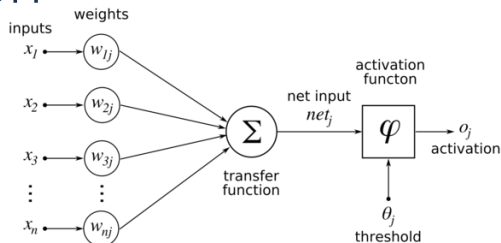


## Notation:

- $x_i$ : the  $i$ 'th input
- $w_i$ : the corresponding synaptic weight
- $\theta$ : the bias (can be eliminated as in the perceptron)

# A model neuron

$$o = \sigma\left(\theta + \sum_p w_p x_p\right),$$



## Notation:

- $x_i$ : the  $i$ 'th input
- $w_i$ : the corresponding synaptic weight
- $\theta$ : the bias (can be eliminated as in the perceptron)
- $o$ : the output

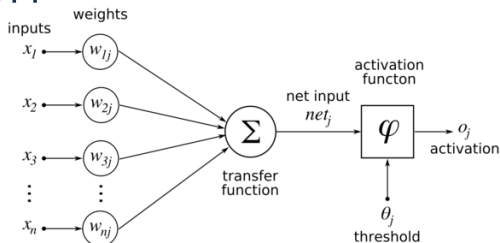
**Learning rule:** Define a loss function  $\ell(\mathbf{y}, \hat{\mathbf{y}})$  for the entire network, and change the weights so as to reduce it. For example, for regression one can use

$$\ell(y, \hat{y}) = (y - \hat{y})^2.$$



# A model neuron

$$o = \sigma\left(\theta + \sum_p w_p x_p\right),$$



## Notation:

- $x_i$ : the  $i$ 'th input
- $w_i$ : the corresponding synaptic weight
- $\theta$ : the bias (can be eliminated as in the perceptron)
- $o$ : the output

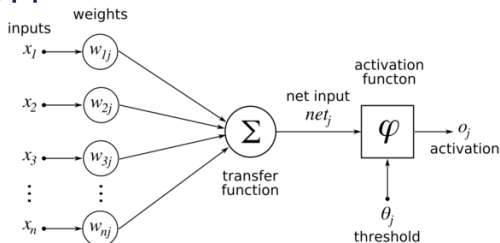
**Learning rule:** Define a loss function  $\ell(\mathbf{y}, \hat{\mathbf{y}})$  for the entire network, and change the weights so as to reduce it. For example, for regression one can use

$$\ell(y, \hat{y}) = (y - \hat{y})^2.$$

The behavior of the network is determined by

# A model neuron

$$o = \sigma\left(\theta + \sum_p w_p x_p\right),$$



## Notation:

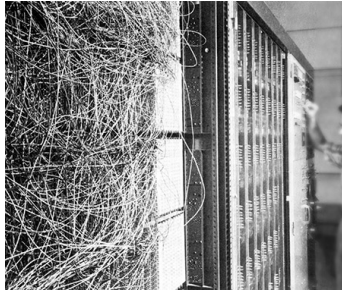
- $x_i$ : the  $i$ 'th input
- $w_i$ : the corresponding synaptic weight
- $\theta$ : the bias (can be eliminated as in the perceptron)
- $o$ : the output

**Learning rule:** Define a loss function  $\ell(\mathbf{y}, \hat{\mathbf{y}})$  for the entire network, and change the weights so as to reduce it. For example, for regression one can use

$$\ell(y, \hat{y}) = (y - \hat{y})^2.$$

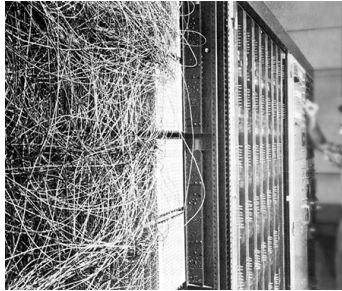
The behavior of the network is determined by

# Multilayer Perceptron



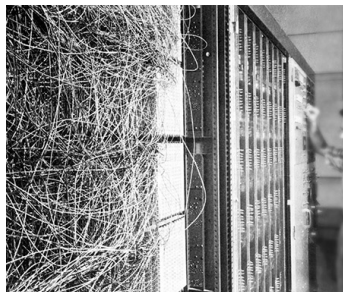
- Multilayer perceptron was thought of as universal model of the brain

# Multilayer Perceptron



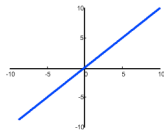
- Multilayer perceptron was thought of as universal model of the brain
- Mark I: 400 pixel image, 2020 photocells

# Multilayer Perceptron



- Multilayer perceptron was thought of as universal model of the brain
- Mark I: 400 pixel image, 2020 photocells
- Problem: hard to train, multilayer perceptron plagued with local minima

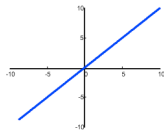
# Activation functions



**Linear:**  $\sigma(t) = t$

Question: What is the problem with this?

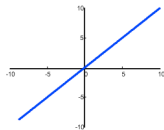
# Activation functions



**Linear:**  $\sigma(t) = t$

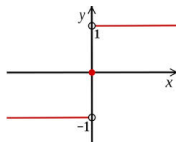
Question: What is the problem with this? Linear functions composed with each other are still linear, so no point in having a multilayer network.

# Activation functions



**Linear:**  $\sigma(t) = t$

Question: What is the problem with this? Linear functions composed with each other are still linear, so no point in having a multilayer network.



**Hard threshold:**  $\sigma(t) = \text{sgn}(t)$

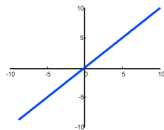
“Threshold Logic Unit” [McCulloch & Pitts, 1943]

→ Perceptron [Rosenblatt, 1958]

Question: What is the problem with this?

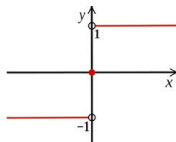


# Activation functions



**Linear:**  $\sigma(t) = t$

Question: What is the problem with this? Linear functions composed with each other are still linear, so no point in having a multilayer network.



**Hard threshold:**  $\sigma(t) = \text{sgn}(t)$

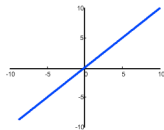
“Threshold Logic Unit” [McCulloch & Pitts, 1943]

→ Perceptron [Rosenblatt, 1958]

Question: What is the problem with this?

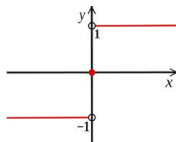
Not differentiable.

# Activation functions



**Linear:**  $\sigma(t) = t$

Question: What is the problem with this? Linear functions composed with each other are still linear, so no point in having a multilayer network.



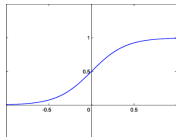
**Hard threshold:**  $\sigma(t) = \text{sgn}(t)$

“Threshold Logic Unit” [McCulloch & Pitts, 1943]

→ Perceptron [Rosenblatt, 1958]

Question: What is the problem with this?

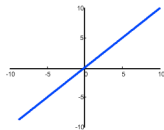
Not differentiable.



**(log-)sigmoid:**  $\sigma(t) = 1/(1 + e^{-t})$

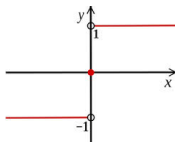
Also called the logistic function.

# Activation functions



**Linear:**  $\sigma(t) = t$

Question: What is the problem with this? Linear functions composed with each other are still linear, so no point in having a multilayer network.



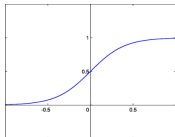
**Hard threshold:**  $\sigma(t) = \text{sgn}(t)$

“Threshold Logic Unit” [McCulloch & Pitts, 1943]

→ Perceptron [Rosenblatt, 1958]

Question: What is the problem with this?

Not differentiable.

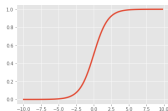


**(log-)sigmoid:**  $\sigma(t) = 1/(1 + e^{-t})$

Also called the logistic function.

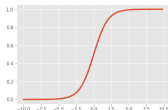
This is what we will use.

# Activation functions

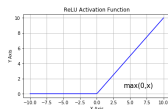


**tanh:**  $\sigma(x) = \tanh(x) = (e^x - e^{-x}) / (e^x + e^{-x})$

# Activation functions

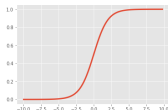


**tanh:**  $\sigma(x) = \tanh(x) = (e^x - e^{-x}) / (e^x + e^{-x})$

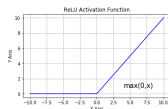


**Rectified linear unit (ReLU):**  $\sigma(x) = \max(0, x)$

# Activation functions



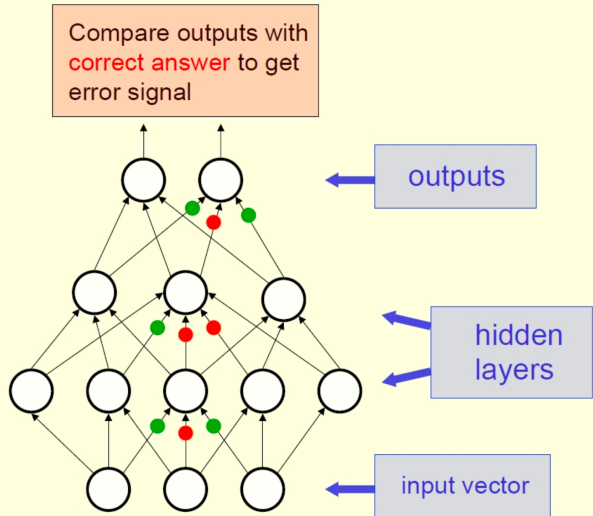
**tanh:**  $\sigma(x) = \tanh(x) = (e^x - e^{-x}) / (e^x + e^{-x})$



**Rectified linear unit (ReLU):**  $\sigma(x) = \max(0, x)$

# Multilayer Neural Nets

Back-propagate  
error signal to  
get derivatives  
for learning



# Multilayer Neural Nets

- Theoretical result: 2-layer net with linear output can approximate any continuous function over compact domain to arbitrary accuracy (given enough hidden units!) [Cybenko 1989]



# Multilayer Neural Nets

- Theoretical result: 2-layer net with linear output can approximate any continuous function over compact domain to arbitrary accuracy (given enough hidden units!) [Cybenko 1989]
- But: there are functions with poly-size logic gate circuit of depth  $k$  that require exponential size when restricted to depth  $O(k)$  [Hastad, 1986]

# Multilayer Neural Nets

- Theoretical result: 2-layer net with linear output can approximate any continuous function over compact domain to arbitrary accuracy (given enough hidden units!) [Cybenko 1989]
- But: there are functions with poly-size logic gate circuit of depth  $k$  that require exponential size when restricted to depth  $O(k)$  [Hastad, 1986]
- In principle, the more layers the better.

# Multilayer Neural Nets

- Theoretical result: 2-layer net with linear output can approximate any continuous function over compact domain to arbitrary accuracy (given enough hidden units!) [Cybenko 1989]
- But: there are functions with poly-size logic gate circuit of depth  $k$  that require exponential size when restricted to depth  $O(k)$  [Hastad, 1986]
- In principle, the more layers the better. But training deep nets is hard  
→ went out of fashion 1990-2006.

# Multilayer Neural Nets

- Theoretical result: 2-layer net with linear output can approximate any continuous function over compact domain to arbitrary accuracy (given enough hidden units!) [Cybenko 1989]
- But: there are functions with poly-size logic gate circuit of depth  $k$  that require exponential size when restricted to depth  $O(k)$  [Hastad, 1986]
- In principle, the more layers the better. But training deep nets is hard  
→ went out of fashion 1990-2006.
- Reasons for resurgence of deep learning:
  - Efficient greedy pretraining idea [Hinton et al., 2006]

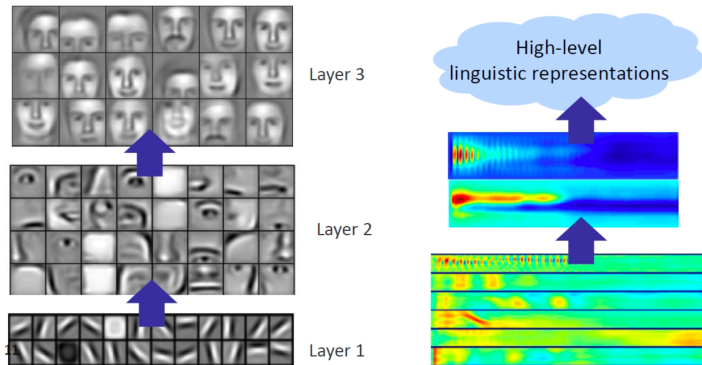
# Multilayer Neural Nets

- Theoretical result: 2-layer net with linear output can approximate any continuous function over compact domain to arbitrary accuracy (given enough hidden units!) [Cybenko 1989]
- But: there are functions with poly-size logic gate circuit of depth  $k$  that require exponential size when restricted to depth  $O(k)$  [Hastad, 1986]
- In principle, the more layers the better. But training deep nets is hard  
→ went out of fashion 1990-2006.
- Reasons for resurgence of deep learning:
  - Efficient greedy pretraining idea [Hinton et al., 2006]
  - Explosion in amount of training data.

# Multilayer Neural Nets

- Theoretical result: 2-layer net with linear output can approximate any continuous function over compact domain to arbitrary accuracy (given enough hidden units!) [Cybenko 1989]
- But: there are functions with poly-size logic gate circuit of depth  $k$  that require exponential size when restricted to depth  $O(k)$  [Hastad, 1986]
- In principle, the more layers the better. But training deep nets is hard  
→ went out of fashion 1990-2006.
- Reasons for resurgence of deep learning:
  - Efficient greedy pretraining idea [Hinton et al., 2006]
  - Explosion in amount of training data.
  - Modern GPU architectures make training much faster.

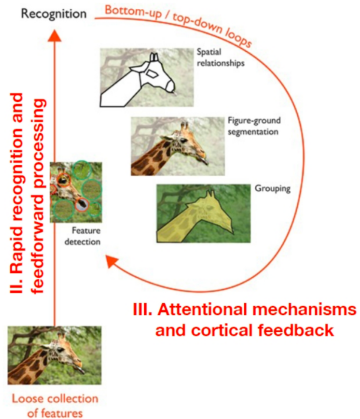
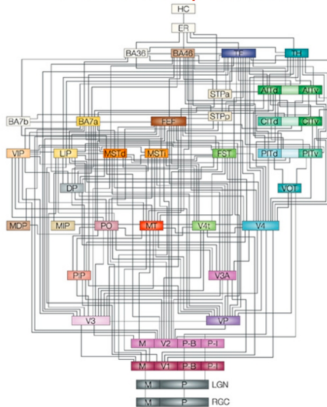
# Multilayer Representations



For vision tasks in particular, representing complex scenes in terms of a hierarchy of features makes sense.

# Multilayer Representations

## Fundamentals of primate vision





# Feed-forward NN for regression

For regression  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ :

- $n$  neurons in input layer (layer 0)

# Feed-forward NN for regression

For regression  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ :

- $n$  neurons in input layer (layer 0)
- $m$  neurons in output layer (layer  $L$ )

# Feed-forward NN for regression

For regression  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ :

- $n$  neurons in input layer (layer 0)
- $m$  neurons in output layer (layer  $L$ )
- Some number of neurons of your choice in the intermediate (hidden) layers

# Feed-forward NN for regression

For regression  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ :

- $n$  neurons in input layer (layer 0)
- $m$  neurons in output layer (layer  $L$ )
- Some number of neurons of your choice in the intermediate (hidden) layers
- Use a loss function like

$$\mathcal{E}(\hat{\mathbf{y}}, \mathbf{y}) = \frac{1}{2} \sum_{i=1}^d (a_{\tau(i)} - y_i)^2,$$

where  $a_{\tau(i)}$  is the output of the  $i$ 'th neuron in the output layer.

# Feed-forward NN for regression

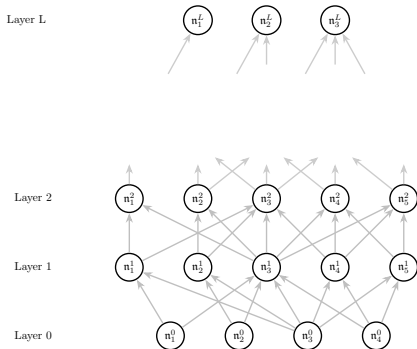
For regression  $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ :

- $n$  neurons in input layer (layer 0)
- $m$  neurons in output layer (layer  $L$ )
- Some number of neurons of your choice in the intermediate (hidden) layers
- Use a loss function like

$$\mathcal{E}(\hat{\mathbf{y}}, \mathbf{y}) = \frac{1}{2} \sum_{i=1}^d (a_{\tau(i)} - y_i)^2,$$

where  $a_{\tau(i)}$  is the output of the  $i$ 'th neuron in the output layer.

- Each neuron may be connected to all the neurons in the previous layer or just a subset of them.



# Feed-forward NN for classification

For regression  $f: \mathbb{R}^n \rightarrow \{1, 2, \dots, k\}$ :

- $n$  neurons in input layer (layer 0)

# Feed-forward NN for classification

For regression  $f: \mathbb{R}^n \rightarrow \{1, 2, \dots, k\}$ :

- $n$  neurons in input layer (layer 0)
- $k$  neurons in output layer (layer  $L$ )

# Feed-forward NN for classification

For regression  $f: \mathbb{R}^n \rightarrow \{1, 2, \dots, k\}$ :

- $n$  neurons in input layer (layer 0)
- $k$  neurons in output layer (layer  $L$ )
- Some number of neurons of your choice in the intermediate (hidden) layers



# Feed-forward NN for classification

For regression  $f: \mathbb{R}^n \rightarrow \{1, 2, \dots, k\}$ :

- $n$  neurons in input layer (layer 0)
- $k$  neurons in output layer (layer  $L$ )
- Some number of neurons of your choice in the intermediate (hidden) layers
- Use a loss function like the log-softmax

$$\mathcal{E}(\hat{\mathbf{y}}, \mathbf{y}) = -\log\left(\frac{e^{a_{\tau(i)}}}{\sum_{i=1}^k e^{a_{\tau(i)}}}\right),$$

where  $a_{\tau(i)}$  is the output of the  $i$ 'th neuron in the output layer.

# Feed-forward NN for classification

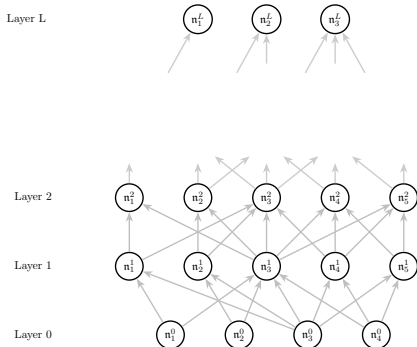
For regression  $f: \mathbb{R}^n \rightarrow \{1, 2, \dots, k\}$ :

- $n$  neurons in input layer (layer 0)
- $k$  neurons in output layer (layer  $L$ )
- Some number of neurons of your choice in the intermediate (hidden) layers
- Use a loss function like the log-softmax

$$\mathcal{E}(\hat{\mathbf{y}}, \mathbf{y}) = -\log\left(\frac{e^{a_{\tau(i)}}}{\sum_{i=1}^k e^{a_{\tau(i)}}}\right),$$

where  $a_{\tau(i)}$  is the output of the  $i$ 'th neuron in the output layer.

- Each neuron may be connected to all the neurons in the previous layer or just a subset of them.



# Training NNs with backpropagation

# Training NNs

- Present training examples one by one (online learning).

# Training NNs

- Present training examples one by one (online learning).
- The error on an example  $(\mathbf{x}, \mathbf{y})$  is some function of the difference between the desired and actual output of the last layer, e.g., for a multivariate regression task

$$\mathcal{E}(\hat{\mathbf{y}}, \mathbf{y}) = \frac{1}{2} \sum_{i=1}^d (a_{\tau(i)} - y_i)^2,$$

where  $\tau(i)$  is the index of the output neuron that is supposed to predict  $y_i$ .

# Training NNs

- Present training examples one by one (online learning).
- The error on an example  $(\mathbf{x}, \mathbf{y})$  is some function of the difference between the desired and actual output of the last layer, e.g., for a multivariate regression task

$$\mathcal{E}(\hat{\mathbf{y}}, \mathbf{y}) = \frac{1}{2} \sum_{i=1}^d (a_{\tau(i)} - y_i)^2,$$

where  $\tau(i)$  is the index of the output neuron that is supposed to predict  $y_i$ .

- Adjust each weight of each neuron in each layer by gradient descent

$$w_{s \rightarrow t} \leftarrow w_{s \rightarrow t} - \eta \frac{\partial \mathcal{E}}{\partial w_{s \rightarrow t}},$$

where  $\eta$  is a parameter called the **learning rate**.

# Notation

Consider a feed-forward architecture with  $L$  layers:

- Set of neurons in layer  $\ell$ :  $\mathcal{N}_\ell$

# Notation

Consider a feed-forward architecture with  $L$  layers:

- Set of neurons in layer  $\ell$ :  $\mathcal{N}_\ell$
- Weight of connection from neuron  $s$  to neuron  $t$ :  $w_{s \rightarrow t}$
- Output (**activation**) of any neuron  $t$  in layer  $\ell$ :

$$a_t = \sigma \left( \underbrace{\sum_{s \in \mathcal{I}(t)} w_{s \rightarrow t} a_s + b_t}_{z_t} \right) = \sigma(z_t),$$

where  $\mathcal{I}_t \subseteq \mathcal{N}_{\ell-1}$  is the set of neurons feeding into  $t$  (in a fully connected feed-forward network  $\mathcal{I}_t = \mathcal{N}_{\ell-1}$ ).



# Notation

Consider a feed-forward architecture with  $L$  layers:

- Set of neurons in layer  $\ell$ :  $\mathcal{N}_\ell$
- Weight of connection from neuron  $s$  to neuron  $t$ :  $w_{s \rightarrow t}$
- Output (**activation**) of any neuron  $t$  in layer  $\ell$ :

$$a_t = \sigma \left( \underbrace{\sum_{s \in \mathcal{I}(t)} w_{s \rightarrow t} a_s + b_t}_{z_t} \right) = \sigma(z_t),$$

where  $\mathcal{I}_t \subseteq \mathcal{N}_{\ell-1}$  is the set of neurons feeding into  $t$  (in a fully connected feed-forward network  $\mathcal{I}_t = \mathcal{N}_{\ell-1}$ ).

- Similarly,  $\mathcal{O}_t$  is the set of neurons in the next layer that  $t$  feeds into.

# Notation

Consider a feed-forward architecture with  $L$  layers:

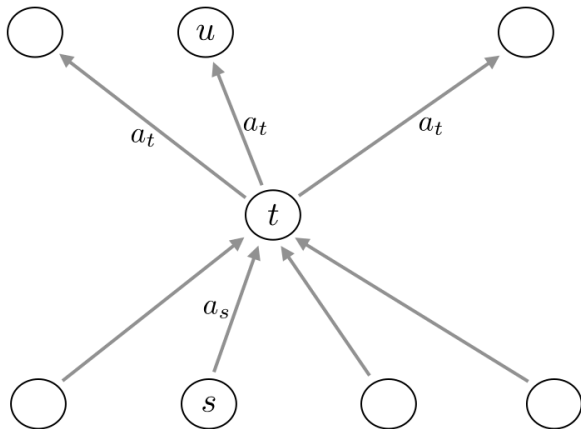
- Set of neurons in layer  $\ell$ :  $\mathcal{N}_\ell$
- Weight of connection from neuron  $s$  to neuron  $t$ :  $w_{s \rightarrow t}$
- Output (**activation**) of any neuron  $t$  in layer  $\ell$ :

$$a_t = \sigma \left( \underbrace{\sum_{s \in \mathcal{I}(t)} w_{s \rightarrow t} a_s + b_t}_{z_t} \right) = \sigma(z_t),$$

where  $\mathcal{I}_t \subseteq \mathcal{N}_{\ell-1}$  is the set of neurons feeding into  $t$  (in a fully connected feed-forward network  $\mathcal{I}_t = \mathcal{N}_{\ell-1}$ ).

- Similarly,  $\mathcal{O}_t$  is the set of neurons in the next layer that  $t$  feeds into.
- The “pre-activation”  $z_t$  plays an important role in the following.

Neurons that  $t$   
feeds into in  
next layer



Neurons in  
previous layer  
that feed into  $t$

# Backpropagation

Problem: how to compute the  $\frac{\partial \mathcal{E}}{\partial w_{s \rightarrow t}}$  derivatives efficiently for all the weights in the neural network?

- Define the **error** (slight misnomer)

$$\delta_t = \frac{\partial \mathcal{E}}{\partial z_t}.$$

# Backpropagation

Problem: how to compute the  $\frac{\partial \mathcal{E}}{\partial w_{s \rightarrow t}}$  derivatives efficiently for all the weights in the neural network?

- Define the **error** (slight misnomer)

$$\delta_t = \frac{\partial \mathcal{E}}{\partial z_t}.$$

- Using the chain rule, we then have the very simple relations

$$\frac{\partial \mathcal{E}}{\partial w_{s \rightarrow t}} = \frac{\partial \mathcal{E}}{\partial z_t} \frac{\partial z_t}{\partial w_{s \rightarrow t}} = \delta_t \cdot a_s, \quad \frac{\partial \mathcal{E}}{\partial b_t} = \delta_t.$$

# Backpropagation

Problem: how to compute the  $\frac{\partial \mathcal{E}}{\partial w_{s \rightarrow t}}$  derivatives efficiently for all the weights in the neural network?

- Define the **error** (slight misnomer)

$$\delta_t = \frac{\partial \mathcal{E}}{\partial z_t}.$$

- Using the chain rule, we then have the very simple relations

$$\frac{\partial \mathcal{E}}{\partial w_{s \rightarrow t}} = \frac{\partial \mathcal{E}}{\partial z_t} \frac{\partial z_t}{\partial w_{s \rightarrow t}} = \delta_t \cdot a_s, \quad \frac{\partial \mathcal{E}}{\partial b_t} = \delta_t.$$

So now the question becomes: how do we compute  $\delta_t$  for all neurons in the network simultaneously?

# Backpropagation

Assume that  $\delta_u$  errors have been computed for all neurons in layer  $\ell + 1$ . Then for any neuron  $t$  in layer  $\ell$ ,

$$\delta_t = \frac{\partial \mathcal{E}}{\partial z_t} = \sum_{u \in \mathcal{O}_t} \underbrace{\frac{\partial \mathcal{E}}{\partial z_u}}_{\delta_u} \frac{\partial z_u}{\partial z_t}.$$

# Backpropagation

Assume that  $\delta_u$  errors have been computed for all neurons in layer  $\ell + 1$ . Then for any neuron  $t$  in layer  $\ell$ ,

$$\delta_t = \frac{\partial \mathcal{E}}{\partial z_t} = \sum_{u \in \mathcal{O}_t} \underbrace{\frac{\partial \mathcal{E}}{\partial z_u}}_{\delta_u} \frac{\partial z_u}{\partial z_t}.$$

and

$$\frac{\partial z_u}{\partial z_t} = \frac{\partial z_u}{\partial a_t} \frac{\partial a_t}{\partial z_t} = w_{t \rightarrow u} \sigma'(z_t).$$



# Backpropagation

Assume that  $\delta_u$  errors have been computed for all neurons in layer  $\ell + 1$ . Then for any neuron  $t$  in layer  $\ell$ ,

$$\delta_t = \frac{\partial \mathcal{E}}{\partial z_t} = \sum_{u \in \mathcal{O}_t} \underbrace{\frac{\partial \mathcal{E}}{\partial z_u}}_{\delta_u} \frac{\partial z_u}{\partial z_t}.$$

and

$$\frac{\partial z_u}{\partial z_t} = \frac{\partial z_u}{\partial a_t} \frac{\partial a_t}{\partial z_t} = w_{t \rightarrow u} \sigma'(z_t).$$

So we have the simple backpropagation formula

$$\delta_t = \sum_{u \in \mathcal{O}_t} w_{t \rightarrow u} \sigma'(z_t) \delta_u.$$

In the last layer,  $\delta_t$  is computed directly from the loss.

# Forward pass

- Initialize the activations of the neurons in layer 0 based on the input  $\mathbf{x}$ .

# Forward pass

- Initialize the activations of the neurons in layer 0 based on the input  $\mathbf{x}$ .
- For layers  $1, 2, \dots, L$  compute the activations from the activations of the previous layer by

$$a_t = \sigma \left( \sum_{s \in \mathcal{I}(t)} w_{s \rightarrow t} a_s + b_t \right) = \sigma(z_t).$$

# Forward pass

- Initialize the activations of the neurons in layer 0 based on the input  $\mathbf{x}$ .
- For layers  $1, 2, \dots, L$  compute the activations from the activations of the previous layer by

$$a_t = \sigma \left( \sum_{s \in \mathcal{I}(t)} w_{s \rightarrow t} a_s + b_t \right) = \sigma(z_t).$$

- Read off  $\hat{\mathbf{y}}$  from the last layer, and compute the loss  $\mathcal{E}(\hat{\mathbf{y}}, \mathbf{y})$ .

# Backward pass

- Initialize the errors in the last layer by directly computing  $\delta_t = \frac{\partial \mathcal{E}}{\partial z_t}$ .

# Backward pass

- Initialize the errors in the last layer by directly computing  $\delta_t = \frac{\partial \mathcal{E}}{\partial z_t}$ .
- For layers  $L - 1, L - 2, \dots, 1$  compute the errors from layer above by

$$\delta_t = \sum_{u \in \mathcal{O}_t} w_{t \rightarrow u} \sigma'(z_t) \delta_u.$$

# Backward pass

- Initialize the errors in the last layer by directly computing  $\delta_t = \frac{\partial \mathcal{E}}{\partial z_t}$ .
- For layers  $L - 1, L - 2, \dots, 1$  compute the errors from layer above by

$$\delta_t = \sum_{u \in \mathcal{O}_t} w_{t \rightarrow u} \sigma'(z_t) \delta_u.$$

- For each neuron in each layer compute the partial derivatives that we ultimately need

$$\frac{\partial \mathcal{E}}{\partial w_{s \rightarrow t}} = \delta_t \cdot a_s, \quad \frac{\partial \mathcal{E}}{\partial b_t} = \delta_t.$$

# Backward pass

- Initialize the errors in the last layer by directly computing  $\delta_t = \frac{\partial \mathcal{E}}{\partial z_t}$ .
- For layers  $L-1, L-2, \dots, 1$  compute the errors from layer above by

$$\delta_t = \sum_{u \in \mathcal{O}_t} w_{t \rightarrow u} \sigma'(z_t) \delta_u.$$

- For each neuron in each layer compute the partial derivatives that we ultimately need

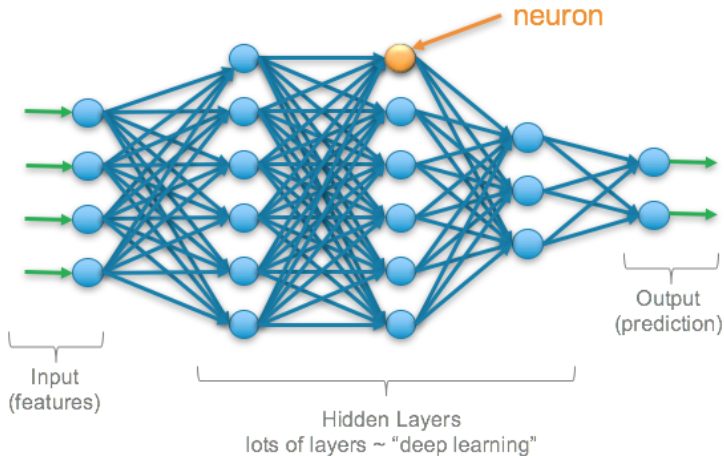
$$\frac{\partial \mathcal{E}}{\partial w_{s \rightarrow t}} = \delta_t \cdot a_s, \quad \frac{\partial \mathcal{E}}{\partial b_t} = \delta_t.$$

- Update the weights and biases by the SGD rule.



# The multi-layer perceptron

---



A feed-forward network where is layer is fully connected is called a multi-layer perceptron.

# Forward pass

For a MLP, the forward iteration can be expressed as

$$\mathbf{a}_\ell = \sigma(\underbrace{W_\ell \mathbf{a}_{\ell-1}}_{\mathbf{z}_\ell} + \mathbf{b}_\ell),$$

where

- $\mathbf{a}_\ell$  is the vector activations in layer  $\ell$ ,
- $\mathbf{b}_\ell$  is the vector of biases in layer  $\ell$ ,
- $W_\ell$  is the matrix of weights from layer  $\ell - 1$  to layer  $\ell$ .

# Forward pass

For a MLP, the forward iteration can be expressed as

$$\mathbf{a}_\ell = \sigma(\underbrace{W_\ell \mathbf{a}_{\ell-1}}_{\mathbf{z}_\ell} + \mathbf{b}_\ell),$$

where

- $\mathbf{a}_\ell$  is the vector activations in layer  $\ell$ ,
- $\mathbf{b}_\ell$  is the vector of biases in layer  $\ell$ ,
- $W_\ell$  is the matrix of weights from layer  $\ell - 1$  to layer  $\ell$ .

The forward pass consists of just computing  $\mathbf{x} \mapsto \mathbf{a}_1 \mapsto \mathbf{a}_2 \mapsto \dots \mapsto \mathbf{a}_L \mapsto \hat{\mathbf{y}} \mapsto \mathcal{E}$ .

# Backward pass

Setting  $\delta_\ell$  as the vector of errors in layer  $\ell$ , the backward iteration is

$$\delta_\ell = \sigma'(z_t) \odot (W_\ell^\top \delta_{\ell+1}),$$

where  $\odot$  is the elementwise product of two vectors.

# Backward pass

Setting  $\delta_\ell$  as the vector of errors in layer  $\ell$ , the backward iteration is

$$\delta_\ell = \sigma'(z_t) \odot (W_\ell^\top \delta_{\ell+1}),$$

where  $\odot$  is the elementwise product of two vectors.

The backward pass consists of computing  $\mathcal{E} \mapsto \delta_L \mapsto \delta_{L-1} \mapsto \dots \mapsto \delta_1$ .

# Backward pass

Setting  $\delta_\ell$  as the vector of errors in layer  $\ell$ , the backward iteration is

$$\delta_\ell = \sigma'(z_t) \odot (W_\ell^\top \delta_{\ell+1}),$$

where  $\odot$  is the elementwise product of two vectors.

The backward pass consists of computing  $\mathcal{E} \mapsto \delta_L \mapsto \delta_{L-1} \mapsto \dots \mapsto \delta_1$ .

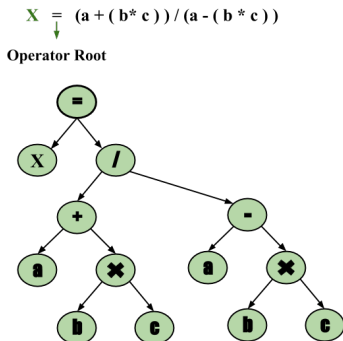
The individual  $\partial\mathcal{E}/\partial w_{s \rightarrow t}$  and  $\partial\mathcal{E}/\partial b_t$  derivatives are computed as before.

# Differentiable computing



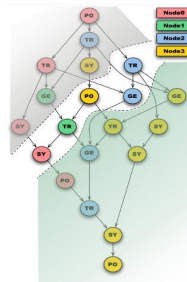
# Computation graphs

It is an old idea in compilers to use a directed acyclic graph (DAG) as an intermediate representation (IR).



Similarly, in high performance computing (HPC) large DAGs are used to depict the interdependencies between parts of a massive compute job.

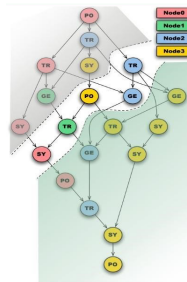
# Computation graphs



The DAG can be used to

- Optimize the order in which computations are performed.

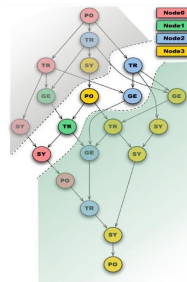
# Computation graphs



The DAG can be used to

- Optimize the order in which computations are performed.
- Optimize the way in which parts of the computation are allocated to processors/nodes.

# Computation graphs



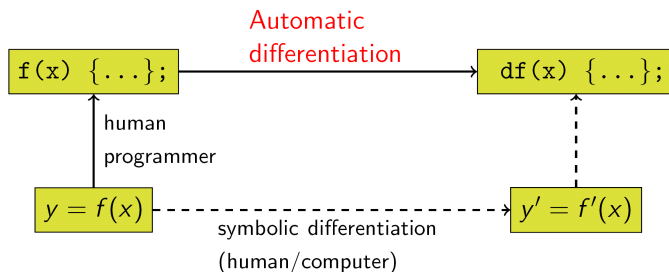
The DAG can be used to

- Optimize the order in which computations are performed.
- Optimize the way in which parts of the computation are allocated to processors/nodes.

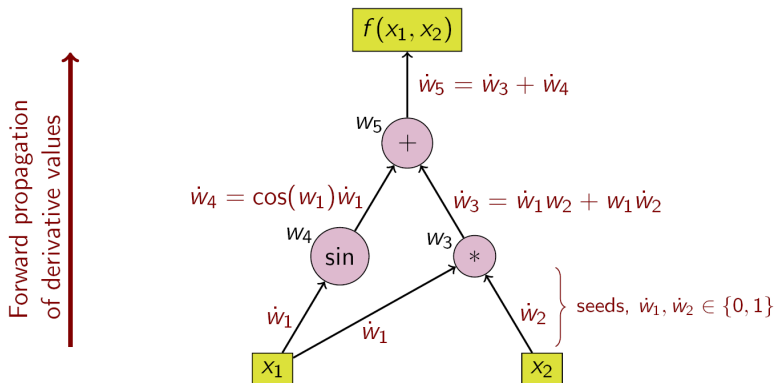
All this can be done statically (ahead of time) or dynamically (at runtime).

# Automatic differentiation

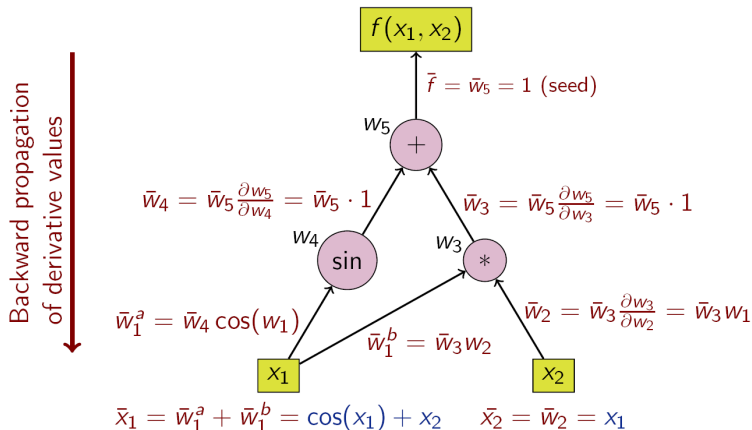
Another interesting field deals with writing compilers that can compute the derivative of any user defined (differentiable) function.



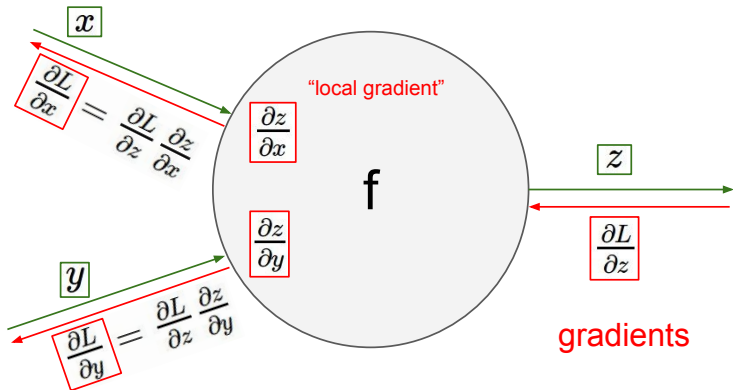
# Automatic differentiation



# Automatic differentiation

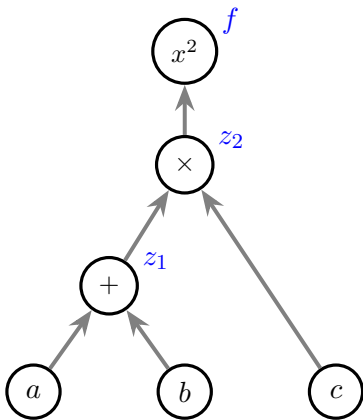


Modern deep learning frameworks combine this idea with DAG-based runtimes.



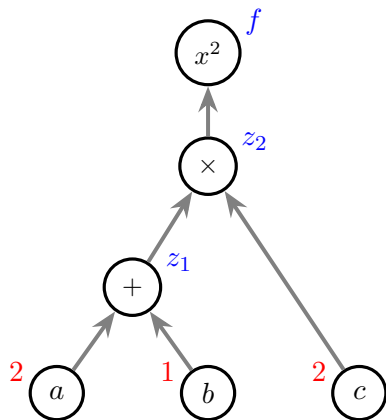


## Forward and backward computations

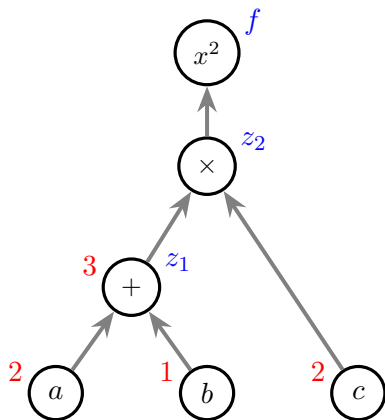


$$f(a, b, c) = ((a + b)c)^2$$

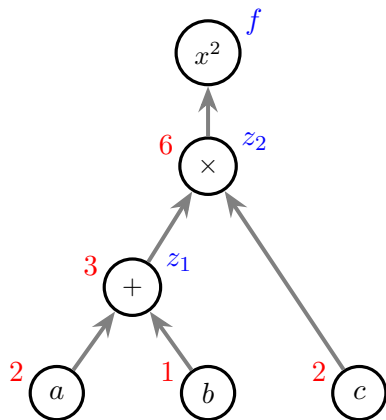
# Forward



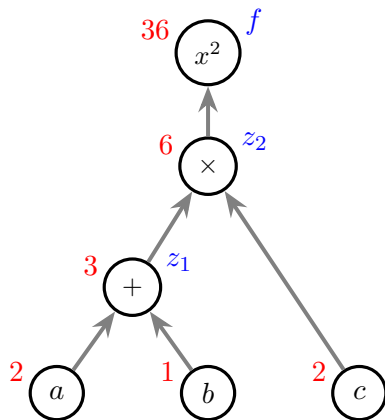
# Forward



# Forward

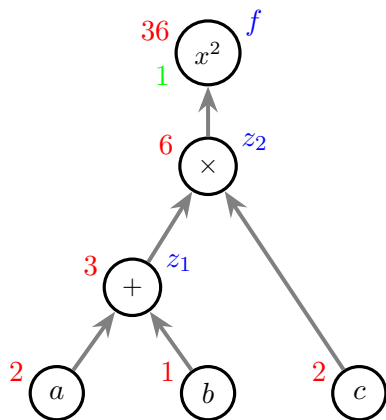


# Forward



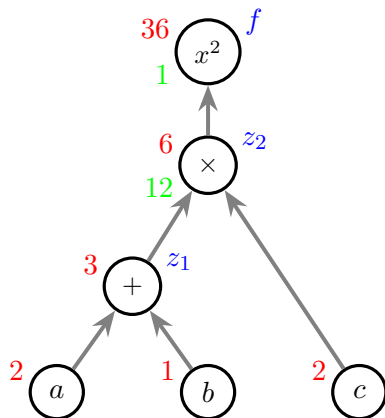
# Backward

$$\frac{\partial f}{\partial f} = 1$$



# Backward

$$f = z_2^2$$

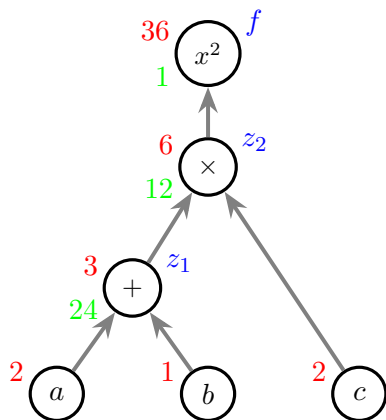


$$\frac{\partial f}{\partial z_2} = \underbrace{\frac{\partial f}{\partial f}}_1 \underbrace{\frac{\partial f}{\partial z_2}}_{2z_2=12} = 12$$



# Backward

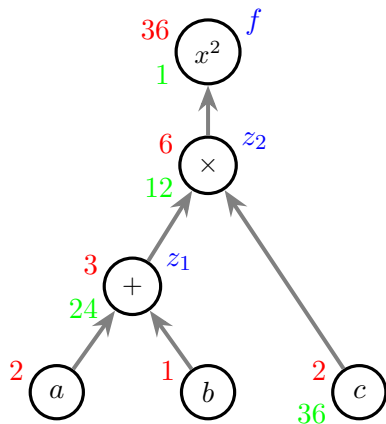
$$z_2 = z_1 \cdot c$$



$$\frac{\partial f}{\partial z_1} = \underbrace{\frac{\partial f}{\partial z_2}}_{12} \underbrace{\frac{\partial z_2}{\partial z_1}}_c = 24$$

# Backward

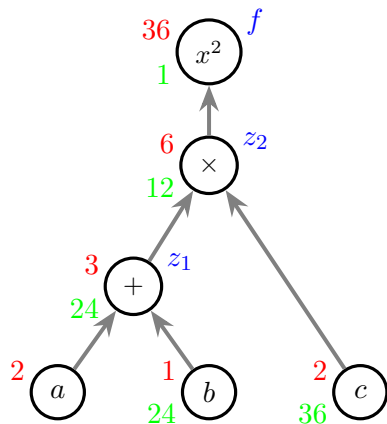
$$z_2 = z_1 \cdot c$$



$$\frac{\partial f}{\partial c} = \underbrace{\frac{\partial f}{\partial z_2}}_{12} \underbrace{\frac{\partial z_2}{\partial c}}_{z_1} = 36$$

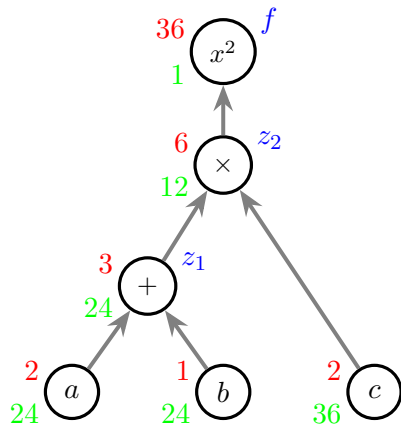
# Backward

$$z_1 = a + b$$



$$\frac{\partial f}{\partial b} = \underbrace{\frac{\partial f}{\partial z_1}}_{24} \underbrace{\frac{\partial z_1}{\partial b}}_1 = 24$$

# Backward



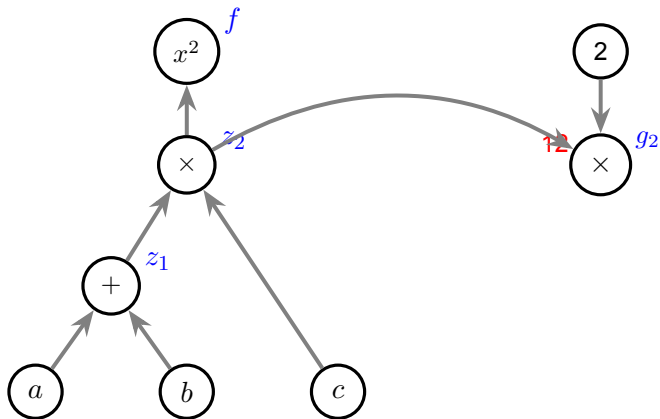
$$z_1 = a + b$$

$$\frac{\partial f}{\partial a} = \underbrace{\frac{\partial f}{\partial z_1}}_{24} \underbrace{\frac{\partial z_1}{\partial a}}_1 = 24$$

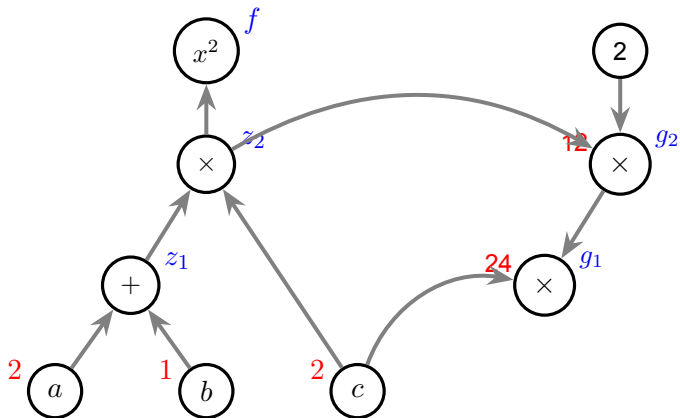
# Symbolic differentiation

---

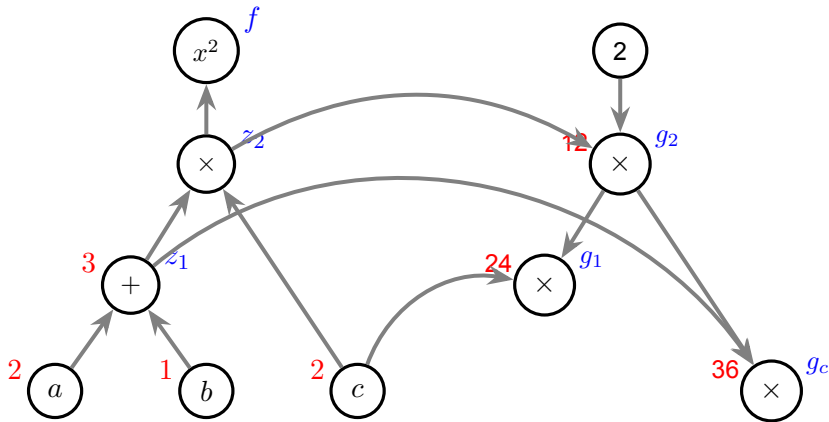
$$g_2 = \frac{\partial f}{\partial z_2} = 2 \cdot z_2$$



$$g_1 = \frac{\partial f}{\partial z_1} = \frac{\partial f}{\partial z_2} \underbrace{\frac{\partial z_2}{\partial z_1}}_c$$

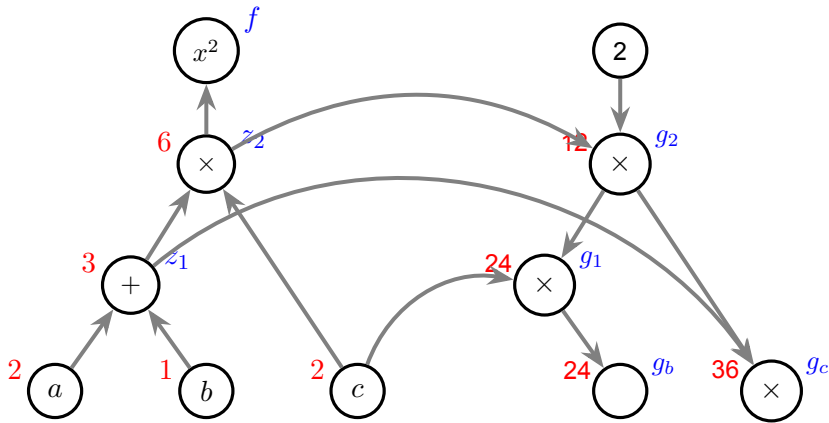


$$g_c = \frac{\partial f}{\partial c} = \frac{\partial f}{\partial z_2} \underbrace{\frac{\partial z_2}{\partial c}}_{z_1}$$

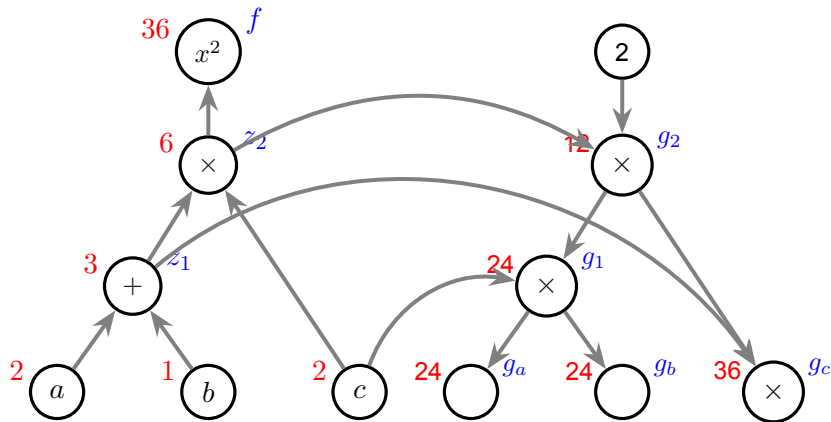




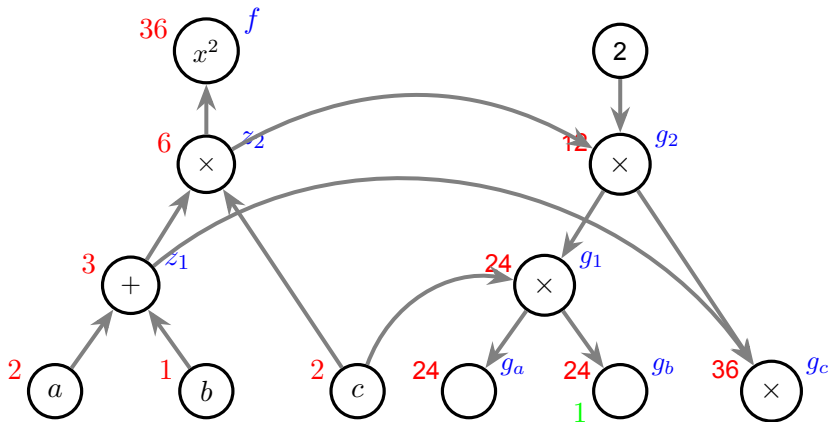
$$g_b = \frac{\partial f}{\partial b} = \frac{\partial f}{\partial z_1} \underbrace{\frac{\partial z_1}{\partial b}}_1$$



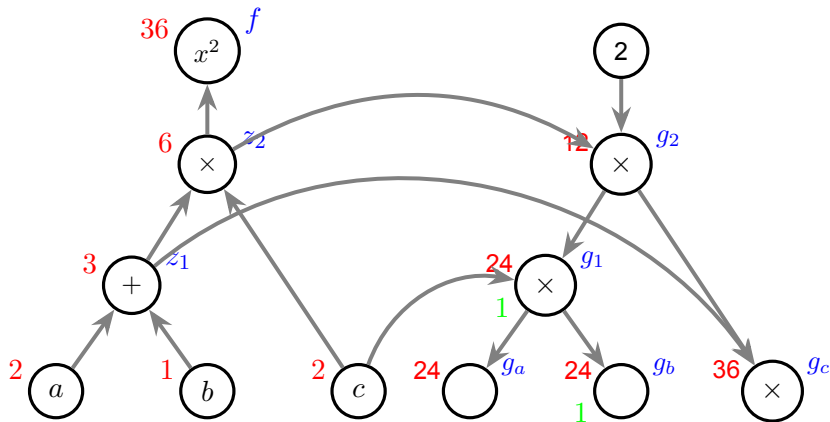
$$g_a = \frac{\partial f}{\partial a} = \frac{\partial f}{\partial z_1} \underbrace{\frac{\partial z_1}{\partial a}}_1$$



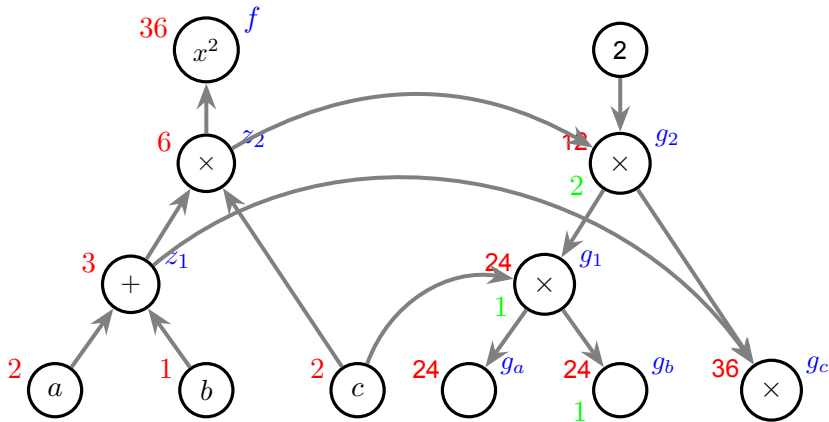
$$\frac{\partial g_b}{\partial g_b} = 1$$



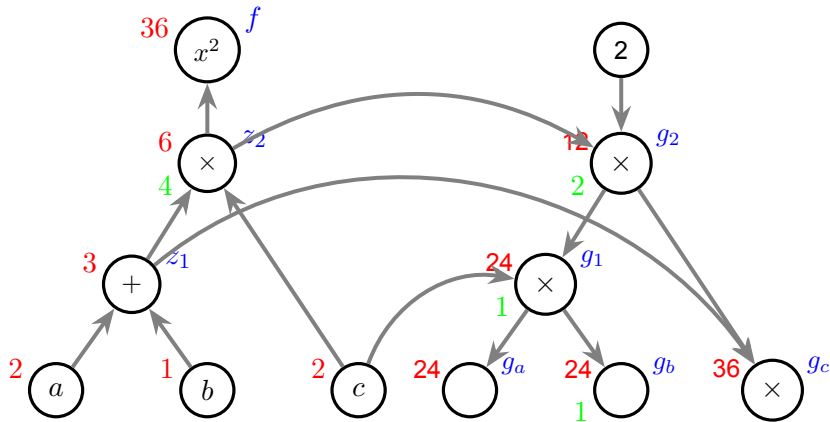
$$\frac{\partial g_b}{\partial g_1} = \frac{\partial g_b}{\partial g_b} = 1$$



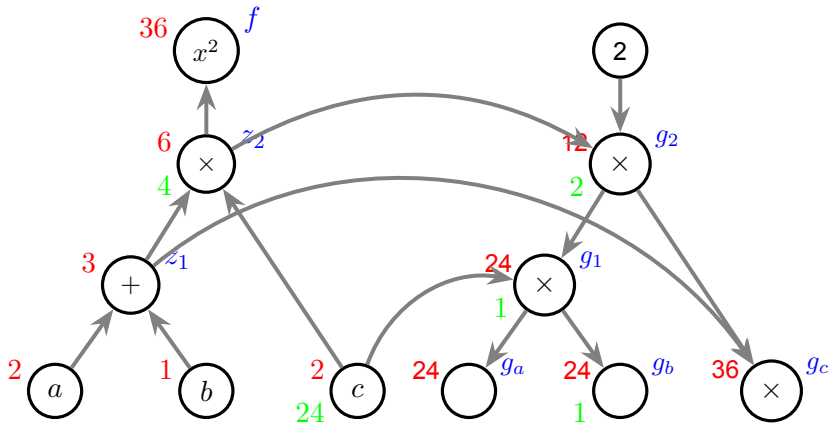
$$\frac{\partial g_b}{\partial g_2} = \underbrace{\frac{\partial g_b}{\partial g_1}}_1 \underbrace{\frac{\partial g_1}{\partial g_2}}_c = 2$$



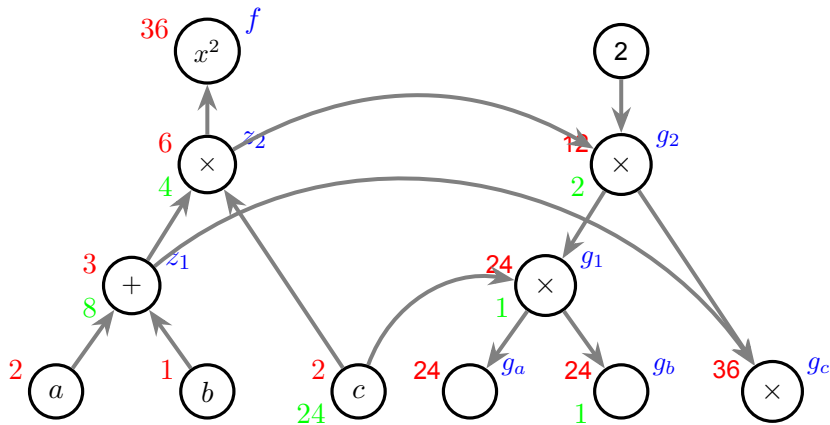
$$\frac{\partial g_b}{\partial z_2} = \underbrace{\frac{\partial g_b}{\partial g_2}}_2 \underbrace{\frac{\partial g_2}{\partial z_2}}_2 = 4$$



$$\frac{\partial g_b}{\partial c} = \underbrace{\frac{\partial g_b}{\partial z_2}}_4 \underbrace{\frac{\partial z_2}{\partial c}}_{z_1} + \underbrace{\frac{\partial g_b}{\partial g_1}}_1 \underbrace{\frac{\partial g_1}{\partial c}}_{g_2} = 4 \cdot 3 + 1 \cdot 12 = 24$$

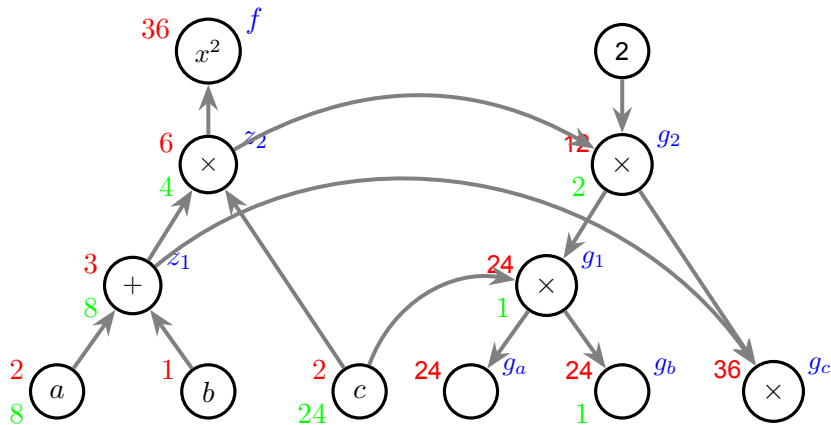


$$\frac{\partial g_b}{\partial z_1} = \underbrace{\frac{\partial g_b}{\partial z_2}}_4 \underbrace{\frac{\partial z_2}{\partial z_1}}_c = 8$$

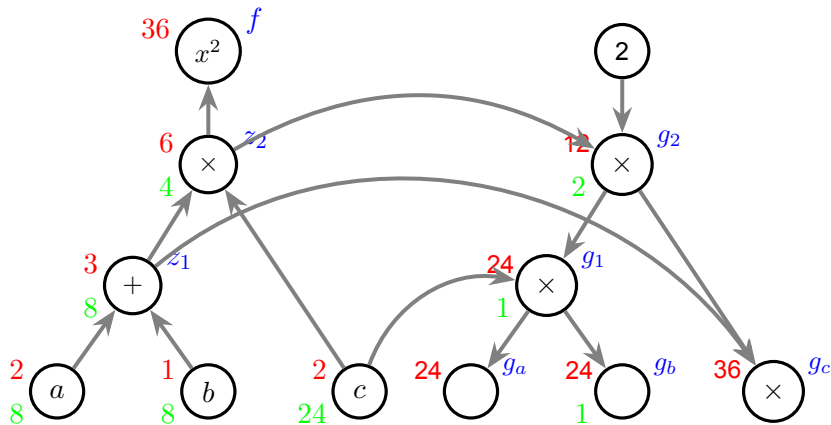




$$\frac{\partial g_b}{\partial a} = \underbrace{\frac{\partial g_b}{\partial z_1}}_8 \underbrace{\frac{\partial z_1}{\partial a}}_1 = 8$$



$$\frac{\partial g_b}{\partial b} = \underbrace{\frac{\partial g_b}{\partial z_1}}_8 \underbrace{\frac{\partial z_1}{\partial b}}_1 = 8$$



# Verification

$$f = ((a + b)c)^2$$

$$\frac{\partial f}{\partial b} = 2(a + b)c^2$$

# Verification

$$f = ((a + b)c)^2$$

$$\frac{\partial f}{\partial b} = 2(a + b)c^2$$

$$\frac{\partial f}{\partial a \partial b} = 2c^2 = 8$$

# Verification

$$f = ((a + b)c)^2$$

$$\frac{\partial f}{\partial b} = 2(a + b)c^2$$

$$\frac{\partial f}{\partial a \partial b} = 2c^2 = 8$$

$$\frac{\partial f}{\partial a \partial c} = 2(a + b)c = 12$$