

06-20416 and 06-12412 (Intro to) Neural Computation

05 – Backpropagation

Per Kristian Lehre

Last lecture

- Functions of multiple variables
- Partial derivatives and the chain rule
- Gradients
 - Direction of steepest ascent
- Gradient descent

This lecture

- Feedforward Neural Networks
 - Model in ML consisting of multiple layers of units
 - Each unit a non-linear transformation of weighted inputs
 - Model parameters are weights and biases
- Backpropagation algorithm (start)
 - How to compute the gradient of the cost of a feedforward neural network wrt its parameters

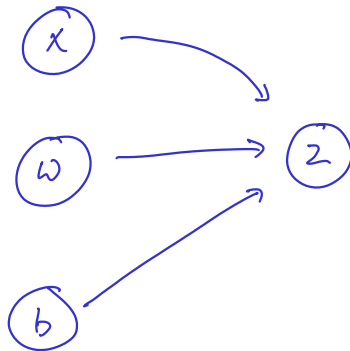
Computation Graphs

We will describe machine learning models using computation graphs where

- nodes represent variables (values, vectors, matrices),
- edges represent functional dependencies, i.e., an edge from x to y indicates that y is a function of x .

Example

A linear regression model $z = \sum_{i=1}^m a_i w_i + b$



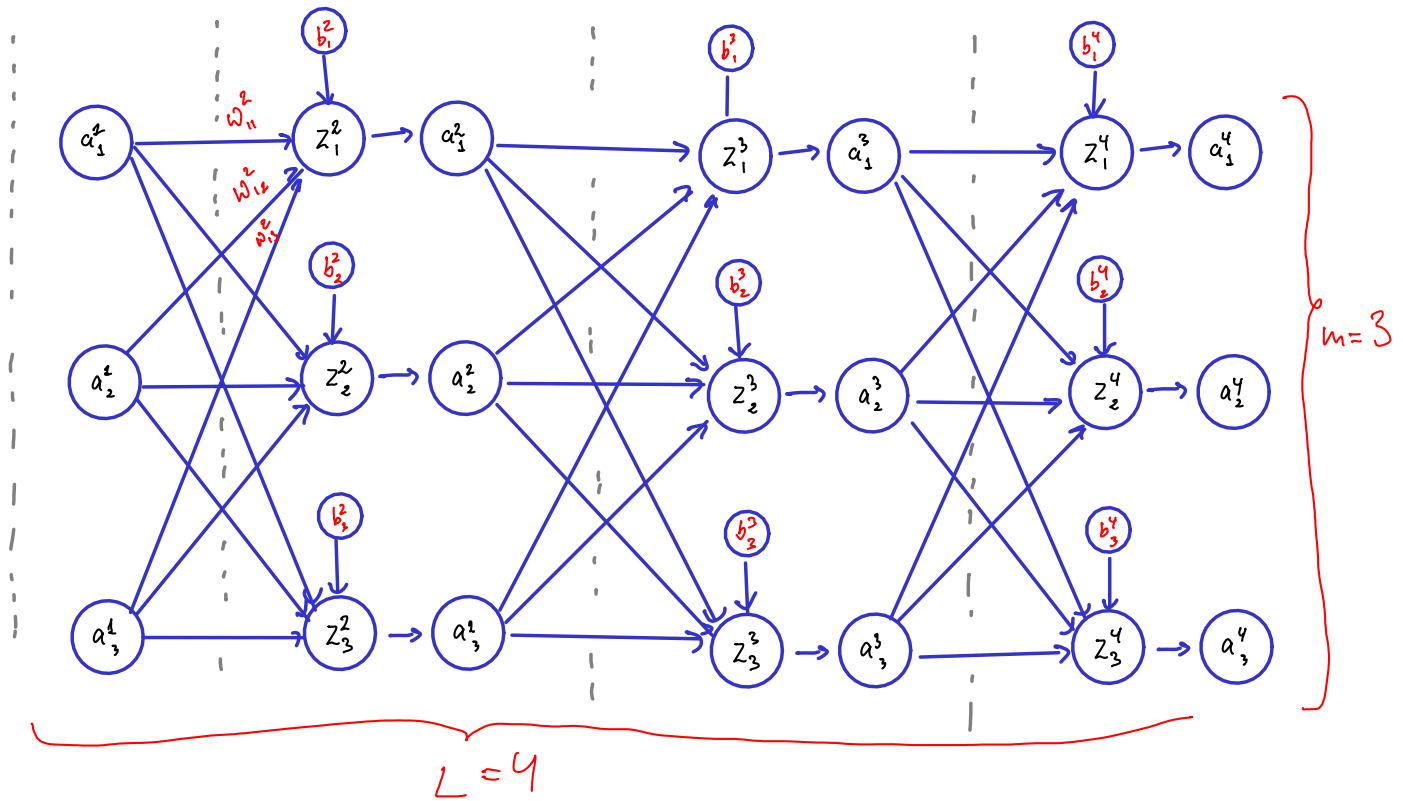
Feed Forward Neural Network

Input
Layer

Hidden
Layer

Hidden
Layer

Output
Layer



L number of layers in the network, where layer 1 is "input layer", and layer L is "output layer"

m "width" of network (can vary between layers)

w_{jk}^l "weight" of connection between k -th unit in layer $l-1$, to j -th unit in layer l

b_j^l "bias" of j -th unit in layer l

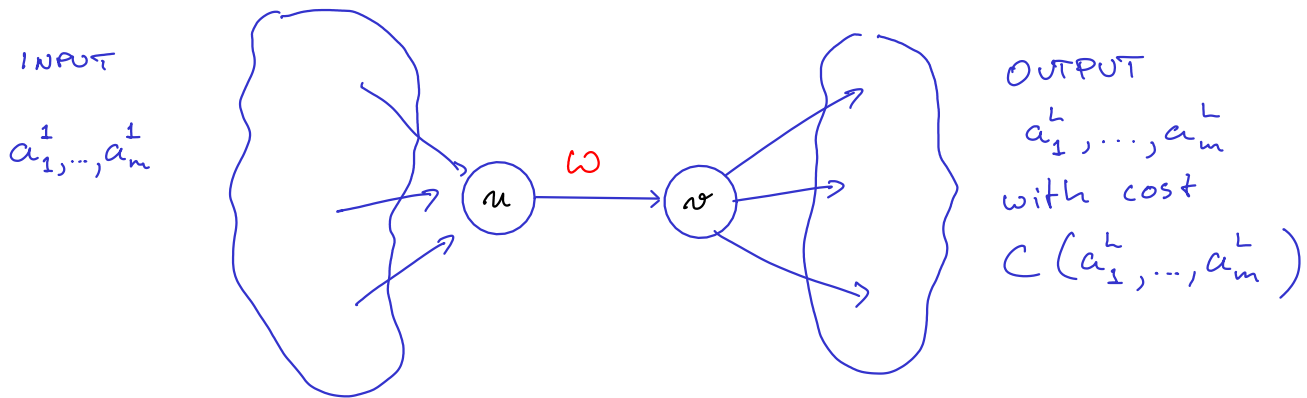
$z_j^l = \sum_k w_{jk}^l a_k^{l-1} + b_j^l$ weighted input to unit j in layer l

$a_j^l = \phi(z_j^l)$ "activation" of unit j in layer l , where ϕ is an "activation function"

Training of Feedforward Neural Networks

The parameters of the network are

- the weights w_{jk}^l in each layer
- the biases b_j^l

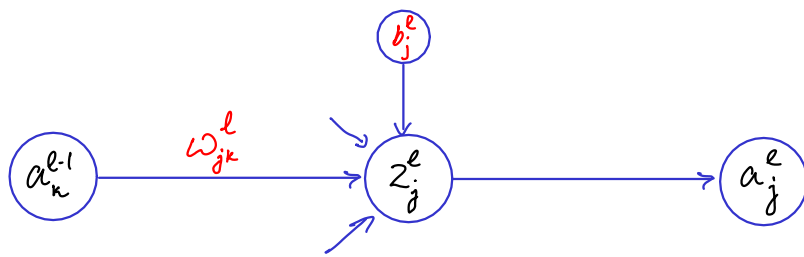


General Idea

To apply gradient descent to optimise a weight w (or bias b) in a network, we apply the chain rule

$$\frac{\partial C}{\partial w} = \frac{\partial C}{\partial v} \cdot \frac{\partial v}{\partial w}$$

Idea applied to Feed Forward Neural Networks



$$z_j^l = \sum_{k=1}^m w_{jk}^l a_k^{l-1} + b_j^l$$

$$a_j^l = \phi(z_j^l)$$

"Activation" of unit k in layer $l-1$.

$$\frac{\partial C}{\partial w_{jk}^l} = \frac{\partial C}{\partial z_j^l} \cdot \frac{\partial z_j^l}{\partial w_{jk}^l} = \frac{\partial C}{\partial z_j^l} \cdot a_k^{l-1}$$

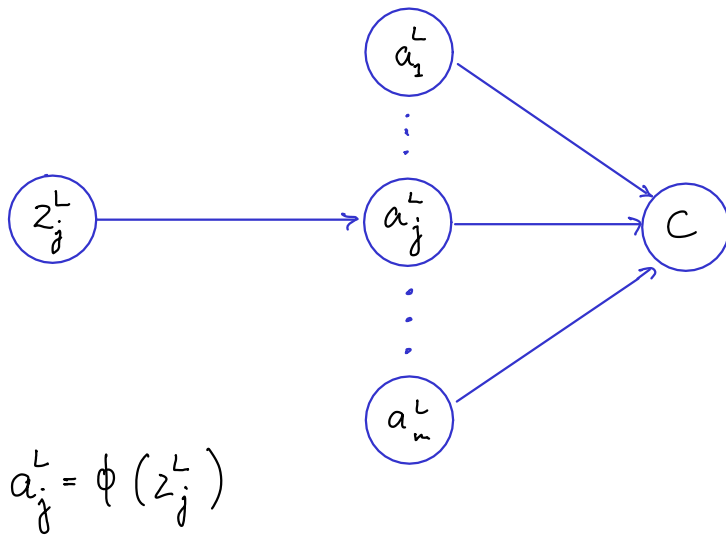
$$\frac{\partial C}{\partial b_j^l} = \frac{\partial C}{\partial z_j^l} \cdot \frac{\partial z_j^l}{\partial b_j^l} = \frac{\partial C}{\partial z_j^l} \cdot 1$$

Hence, we can compute $\frac{\partial C}{\partial w_{jk}^l}$ and $\frac{\partial C}{\partial b_j^l}$ if we know

$$\delta_j^l := \frac{\partial C}{\partial z_j^l}$$

The vector δ^l is called the local gradient for layer l .

Local Gradient for Output Layer



The local gradient for the output layer is

$$\begin{aligned}\delta_j^L &= \frac{\partial C}{\partial z_j^L} && \text{by definition} \\ &= \frac{\partial C}{\partial a_j^L} \cdot \frac{\partial a_j^L}{\partial z_j^L} && \text{by the chain rule} \\ &= \frac{\partial C}{\partial a_j^L} \cdot \phi'(z_j^L) && \text{because } a_j^L = \phi(z_j^L)\end{aligned}$$

The partial derivative $\frac{\partial C}{\partial a_j^L}$ depends on the cost function.

For example, for a regression problem in m dimensions, one could define

$$C(a_1^L, \dots, a_m^L) := \frac{1}{2} \sum_{k=1}^m (y_k - a_k^L)^2,$$

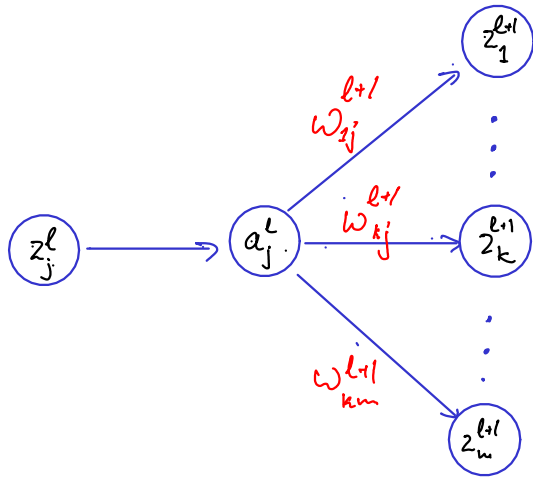
in which case

$$\frac{\partial C}{\partial a_j^L} = a_j^L - y_j$$

desired
output in
 k -th dimension

predicted
output
in k -th
dimension

Local Gradient for Hidden Layers



$$z_k^{l+1} = \sum_r w_{kr}^{l+1} a_r^l$$

$$a_j^l = \phi(z_j^l)$$

$$\delta_j^l = \frac{\partial C}{\partial z_j^l}$$

$$= \frac{\partial C}{\partial a_j^l} \cdot \frac{\partial a_j^l}{\partial z_j^l}$$

$$= \left(\sum_k \frac{\partial C}{\partial z_k^{l+1}} \cdot \frac{\partial z_k^{l+1}}{\partial a_j^l} \right) \cdot \phi'(z_j^l)$$

$$= \phi'(z_j^l) \sum_k \delta_k^{l+1} \cdot w_{kj}^{l+1}$$

by def of local gradient δ_j^l

by chain rule

by chain rule w.r.t $\frac{\partial C}{\partial a_j^l}$

by def. of local gradient δ_k^{l+1}

Summary

For all weights and biases

$$\frac{\partial C}{\partial \omega_{jk}^l} = \delta_j^l \cdot a_k^{l-1}$$

$$\frac{\partial C}{\partial b_j^l} = \delta_j^l$$

where the local gradient δ_j^l is

$$\delta_j^l = \begin{cases} \phi'(z_j^L) \cdot \frac{\partial C}{\partial a_j^L} & \text{if } l=L \text{ (output layer)} \\ \phi'(z_j^l) \cdot \sum_k \delta_k^{l+1} \cdot \omega_{kj}^{l+1} & \text{otherwise (hidden layer)} \end{cases}$$

Matrix Description

The backpropagation algorithm can exploit efficient implementations of matrix operations, e.g. in numerical libraries such as NumPy or on a GPU.

Recall that for a matrix $A \in \mathbb{R}^m \times \mathbb{R}^n$, A_{ij} denotes the element in the i -th row and j -th column.

For two matrices $A \in \mathbb{R}^m \times \mathbb{R}^n$ and $B \in \mathbb{R}^n \times \mathbb{R}^l$

$$(A^T)_{ij} = A_{ji}$$

matrix transpose

$$(AB)_{ij} = \sum_k A_{ik} B_{kj}$$

matrix multiplication

For two vectors $u, v \in \mathbb{R}^m$

$$u + v = (u_1 + v_1, \dots, u_m + v_m)$$

vector addition

$$u \cdot v = \sum_{i=1}^m u_i v_i$$

dot product

$$u \odot v = (u_1 v_1, \dots, u_m v_m)$$

Hadamard product

$$(uv^T)_{ij} = u_i v_j$$

outer product

Weighted Inputs and Activations

$$z^l = (z_1^l, \dots, z_m^l)$$

$$= \left(\sum_{k=1}^m w_{1k}^l a_k^{l-1} + b_1^l, \dots, \sum_{k=1}^m w_{mk}^l a_k^{l-1} + b_m^l \right)$$

$$= w^l a^{l-1} + b$$

$$a^l = (a_1^l, \dots, a_m^l)$$

$$= (\phi(z_1^l), \dots, \phi(z_m^l))$$

$$= \phi(z^l)$$

Local Gradients

Output Layer

$$\begin{aligned}\delta^L &= (\delta_1^L, \dots, \delta_m^L) \\ &= \left(\frac{\partial C}{\partial a_1^L} \cdot \phi'(z_1^L), \dots, \frac{\partial C}{\partial a_m^L} \cdot \phi'(z_m^L) \right) \\ &= \nabla_{a^L} C \odot \phi'(z^L)\end{aligned}$$

Hidden Layer

$$\begin{aligned}\delta^l &= (\delta_1^l, \dots, \delta_m^l) \\ &= \left(\phi'(z_1^l) \cdot \sum_k \delta_k^{l+1} \cdot w_{k1}^{l+1}, \dots, \phi'(z_m^l) \cdot \sum_k \delta_k^{l+1} \cdot w_{km}^{l+1} \right) \\ &= \phi'(z^l) \odot \left(\sum_k (w^{l+1})_{1k}^T \delta_k^{l+1}, \dots, \sum_k (w^{l+1})_{mk}^T \delta_k^{l+1} \right) \\ &= \phi'(z^l) \odot (w^{l+1})^T \delta^{l+1}\end{aligned}$$

Backpropagation Algorithm

Input: A training example $(x, y) \in \mathbb{R}^m \times \mathbb{R}^{m'}$

1. Set the activation in the input layer

$$a^1 = x$$

2. for each $l=2$ to L , feed forward

$$z^l = w^l a^{l-1} + b^l$$

$$a^l = \phi(z^l)$$

3. compute local gradient for output layer

$$\delta^L := \nabla_a^L C \odot \phi'(z^L)$$

4. backpropagate local gradients for hidden layers, i.e.

for each $l=L-1$ to 2

$$\delta^l := ((w^{l+1})^T \delta^{l+1}) \odot \phi'(z^l)$$

5. return the partial derivatives

$$\frac{\partial C}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l$$

$$\frac{\partial C}{\partial b_j^l} = \delta_j^l$$

Training Feed Forward Networks

Assume n training examples

$$(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)}),$$

and a cost function

$$C = \frac{1}{n} \sum_{i=1}^n C^{(i)},$$

where $C^{(i)}$ is the cost on the i -th example.

For example, with MSE, we can define

$$C^{(i)} = \frac{1}{2} (y^{(i)} - a^L)^2$$

where a^L is the output of the network when $x^{(i)}$

Backpropagation gives us the gradient of the overall cost function as follows

$$\frac{\partial C}{\partial w^e} = \frac{1}{n} \sum_{i=1}^n \frac{\partial C^{(i)}}{\partial w^e}$$

$$\frac{\partial C}{\partial b^e} = \frac{1}{n} \sum_{i=1}^n \frac{\partial C^{(i)}}{\partial b^e}$$

NB "Average" gradient
per training example.

In principle, we can now use gradient descent to optimise the weights w and biases b .

Mini-batch Gradient Descent

Computing the gradients is expensive when the number of training examples n is large.

We can approximate the gradients

$$\frac{\partial C}{\partial w^l} \approx \frac{1}{b} \sum_{i=1}^b \frac{\partial C^{(i)}}{\partial w^l}$$

$$\frac{\partial C}{\partial b^l} \approx \frac{1}{b} \sum_{i=1}^b \frac{\partial C^{(i)}}{\partial b^l}$$

using a random "mini-batch" of $b \leq n$ training examples.

$1 < b < n \Rightarrow$ Mini-batch Gradient Descent

$b = 1 \Rightarrow$ Stochastic Gradient Descent

$b = n \Rightarrow$ Batch Gradient Descent

Common to use mini-batch size $b \in (20, 100)$.

Summary

- Computation graphs
- Feedforward Neural Networks
 - Input nodes and biases
 - Activation functions
 - Input, hidden, and output layers
- Backpropagation algorithm (start)
 - The gradient can be computed using local derivatives.
 - Local derivatives are computed backwards, starting from the output layer

Next lecture

- More about backpropagation